

# Chapter 5: Data Science

## Table of Contents

1		
2	<b>Table of Contents</b>	
3	Introduction	3
4	<b>The Statistical and Data Science Investigation Process</b>	8
5	Vignette 1: CODAP	10
6	Vignette 2: Dear Data	13
7	Data Talks K–12	17
8	Transitioning from Pre–K	21
9	K–5	21
10	What questions can data help to answer?	22
11	Asking Questions, Collecting and Analyzing data	25
12	Interpreting and Communicating Results	26
13	Preparing for the major data science work of grades 6–8	28
14	Vignette: Logan from Kindergarten through Grade 5	30
15	Grades 6–8	34
16	Data in the world: Question asking, exploration, interpretation, decision making,	
17	ethics, technology	35
18	Describing, displaying, and comparing statistical variability (grades 6–7)	36
19	Sampling to understand a population: randomness, bias, how many? (grades 7–8)	
20		39
21	Are they related? Two changing quantities (grade 8)	42
22	What are the chances? Probability as the basis for data-based claims	43
23	Vignette	45

24	High School	46
25	Data science for equity and inclusion	47
26	Data for all: living in an information-overloaded world	50
27	Interpreting Categorical and Quantitative Data	52
28	Making Inferences and Justifying Conclusions	55
29	From Statistics to Data Science	57
30	Advanced high school data science	58
31	Design Principles	59
32	Content Learning Outcomes	68
33	Understanding the Role of Data in The World	68
34	Asking Data-Based Questions	69
35	Unraveling the Story That Data Is Telling	69
36	Grappling with Variability and Uncertainty	71
37	Transforming Data with Technology	72
38	Sample Courses	72
39	High School Tools and Resources	74
40	Conclusion	76
41	Free Resources for the Teaching of Data Science	77
42	References	77

43

44 [Note from writers to CFCC for December 16–17 meeting: This chapter still needs major  
45 revision, especially in the high school section. Two California members of the writing  
46 team for the American Statistical Association’s brand-new report “Pre-K–12 Guidelines  
47 for Assessment and Instruction in Statistics Education II (GAISE II): A Framework for  
48 Statistics and Data Science Education” have assisted the writing team by identifying the  
49 major changes needed to the chapter to bring it more in line with modern data  
50 science/statistics education research as reflected in GAISE II. The “statistical problem-  
51 solving process” has been partially integrated into this revision; the many roles of

52 questioning have been made clearer and more prominent; and the definition of data  
53 science has been improved. More work needs to be done, to make sure the chapter is  
54 as accurate as possible in setting out the K–12 needs for this emerging field. In addition,  
55 many CFCC changes are waiting on completion of the more structural changes.]

## 56 Introduction

57 The ability to work with and understand data has become an essential life skill in our  
58 newly data-filled world. Students participate in a world driven by data; making sense of  
59 data, being able to identify data that is misleading, and using data to make decisions  
60 are all important aspects of their role as global citizens in the larger world. It is not only  
61 those who have careers in data science—almost all occupations now require that  
62 employees collect feedback from data and adjust their practice. Stories about the world  
63 are illuminated by massive quantities of data, and community members telling and  
64 listening to those stories need to be able to make sense of data to understand their  
65 health, finances, and news feeds.

66 The numbers are staggering: around 1.7 megabytes of digital data were created and  
67 stored *every second for every person on earth* in 2020, and the vast majority of data  
68 goes unanalyzed (<https://techjury.net/stats-about/big-data-statistics/>). Our lives are  
69 increasingly subject to data-driven algorithms that determine much about our daily  
70 experience, including what ads we see, which neighborhoods receive business or public  
71 investment, who gets screened more closely at the airport, who receives favorable loan  
72 terms, and which medical procedures are recommended or approved.

73 All California students should graduate from high school with data literacy and options  
74 to learn an introduction to data science. Data literacy refers to the ability to reason with  
75 and about data, to make good decisions based on data, to ask questions of data, and to  
76 use statistical reasoning. Data science is an emerging discipline that includes  
77 understanding principles of data collection, data manipulation, data analysis, inference,  
78 and interpretation and communication. The Common Core State Standards set out the  
79 learning of statistics K–12. A data science lens can help the statistical ideas in the  
80 standards come alive and have relevance and meaning for students.

81 GAISE II is a professional report from the American Statistical Association (ASA) setting  
82 out guidelines for assessment and instruction K–12 in statistics and data science, and is  
83 an important resource for this area of mathematics. In the GAISE II report they  
84 emphasize the following:

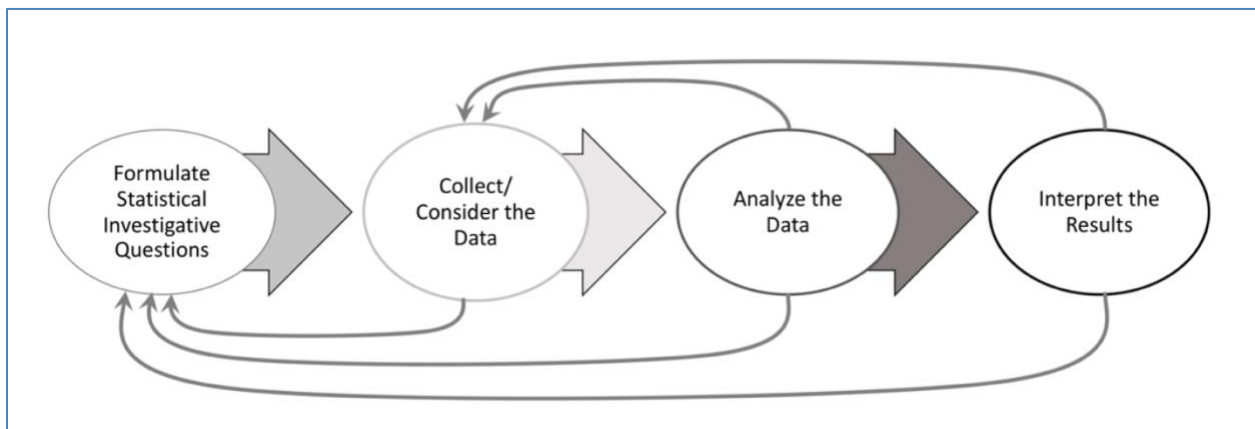
- 85 1. The importance of asking questions throughout the statistical problem-solving  
86 process (formulating a statistical investigative question, collecting or considering  
87 data, analyzing data, and interpreting results), and how this process remains at  
88 the forefront of statistical reasoning for all studies involving data
- 89 2. The consideration of different data and variable types, the importance of carefully  
90 planning how to collect data or how to consider data to help answer statistical  
91 investigative questions, and the process of collecting, cleaning, interrogating, and  
92 analyzing the data
- 93 3. The inclusion of multivariate thinking throughout all Pre-K–12 educational levels
- 94 4. The role of probabilistic thinking in quantifying randomness throughout all levels
- 95 5. The recognition that modern statistical practice is intertwined with technology,  
96 and the importance of incorporating technology as feasible
- 97 6. The enhanced importance of clearly and accurately communicating statistical  
98 information
- 99 7. The role of assessment at the school level, especially items that measure  
100 conceptual understanding and require statistical reasoning involving the  
101 statistical problem-solving process. (GAISE II, 2020, p. 2)

102 Students should be able to draw on the Standards for Mathematical Practices (SMP)  
103 through a statistical lens articulated in the ASA’s Statistical Education of Teachers  
104 (SET) report. For example, students should reason abstractly and quantitatively by  
105 engaging in statistical thinking while considering where data come from (SMP.2), apply  
106 statistical models to “include descriptions of the variability present in data (SMP.4), and  
107 consider available tools such as calculators, spreadsheets, applets, statistical  
108 packages, and graphical displays to help facilitate the statistical problem-solving  
109 process (SMP.5). When students participate in the analysis of large datasets, they  
110 should be able to decide which questions matter, and identify which ones can be

111 answered with a given dataset (SMP.4). The statistical problem-solving process is used  
112 within the process. Further, students should understand some of the ways in which data  
113 are frequently misunderstood or misused and should understand the content and  
114 implications of their own digital data footprints. Finally, students should be prepared to  
115 pursue additional study directed towards fields which include more intensive work with  
116 data, such as designing data collection, deciding on statistical measures appropriate to  
117 the questions under consideration, or making conclusions and claims based on data.

118 The statistical and data science problem solving process, as set out in GAISE II, is  
119 shown in Figure 1:

120 Figure 1. The statistical and data science problem solving process, GAISE II



121

122 The California Common Core State Standards in Mathematics (CA CCSSM) contain  
123 many content standards that help to build the data understanding and skills that high  
124 school graduates require. However, the progression—from counting, categorizing, and  
125 simple picture graphs, to the complex skills and understanding that older students may  
126 develop—requires careful thought and considerably more focus through the K–12  
127 curriculum than most students have historically experienced. The study of data today is  
128 broader than it has ever been. The types of data being collected are vast and the types  
129 of techniques used to analyze data can now rely heavily on computational tools. The  
130 statistical problem-solving process is important as it provides the foundation for finding  
131 meaning in data. Data science and statistics are the science of working with data. The

132 development of statistics and data science mastery articulated in this chapter  
133 represents a modern lens through which to examine the CA CCSSM.

134 Educators regularly use data at the student and classroom level to try to drive  
135 instructional decisions. However, a data science perspective can help educators create  
136 experiences in which their students learn to “read and write the world with mathematics”  
137 (Gutstein 2003). As emphasized throughout this framework, students must experience  
138 mathematics as tools for making sense of and impacting their worlds.

139 Educators should be encouraged to bring data science and statistics directly into their  
140 classroom to create student experiences that are meaningful. Students can experience  
141 statistics and data science as tools for making sense of and impacting their worlds. The  
142 statistical problem-solving process (GAISE II) helps students formulate statistical  
143 investigative questions, take in information by collecting primary data or considering  
144 secondary data, analyze the data to identify relationships and patterns, and (in many  
145 cases) interpret results to answer the question and propose changes to impact the way  
146 the world works.

147 Students who are exposed to and have the capacity to understand data concepts at an  
148 early age begin to develop data literacy and data sense in parallel with number sense.  
149 As students progress through school they should learn different approaches to data  
150 analysis culminating in the investigation of large data sets using up-to-date  
151 technological tools.

152 As students learn the investigative statistical and data science process they should  
153 always consider meaning and context. In the past, some learning of statistics was  
154 removed from situational settings, leading students to learn abstract methods. Data  
155 science involves developing meaning and communicating about a data-rich situation; it  
156 should never be removed from its context. Teachers can use local data sets that give  
157 students the opportunity to ask questions that are meaningful to them, that can help  
158 their local community, or school, allowing students to experience using mathematics to  
159 be an engaged citizen. Statistics and data science is about studying situations—asking  
160 questions such as: Who collected the data? How was it collected? What is the unit of

161 analysis? Teachers can ask students to turn and talk to their partners and groups about  
162 these questions.

163 In this chapter, we present the progression of data literacy and data science standards  
164 and the types of experiences that help build the necessary skills and understandings.  
165 Four important principles in the learning of data science are these:

- 166 1. Students should experience working with data from a context that is meaningful  
167 to them personally. They should have opportunities to solve problems of value to  
168 the students and to their schools and communities.
- 169 2. Students should learn to engage with real data that include multiple variables. At  
170 first students can learn to understand two variables with bivariate data, as they  
171 progress through the grades they can learn to handle multivariable data and  
172 multivariate thinking.
- 173 3. Data investigations should be investigative and collaborative, with students  
174 working together to learn the data science and statistical investigative process.
- 175 4. Familiarity with technology and modern tools should progress through the  
176 grades.

177 As discussed in more detail in the Chapter 2: Teaching for Equity and Engagement, it is  
178 more effective for teachers to plan around big ideas than sets of mathematical methods,  
179 and to choose rich tasks that elicit big ideas. In this chapter we set out the big ideas of  
180 data science that build to the kind of connected understanding needed.

181 **Definition:** Data are observations or measurements in context. **Usage note:** In Latin,  
182 the word *data* is the plural of *datum*. However, in English, *data* is now also commonly  
183 used with singular verbs and refers to a collection of data points. Thus, “the data shows  
184 a correlation...” is more common than “the data show a correlation....” In this chapter  
185 we most often use the word *data* in this way—to refer to a collection of data points—and  
186 in these contexts it takes singular verbs.

187 **Sources:** The development of data science described and illustrated here is guided by  
188 the California Common Core State Standards in Mathematics, and is also informed by  
189 and largely consistent with the following documents:

- 190 ● The Guidelines for Assessment and Instruction in Statistics Education Pre-K–12
- 191 Report (Bargagliotti, Franklin, Arnold, Gould, Johnson, Perez, & Spangler 2020;
- 192 [https://www.amstat.org/asa/education/Guidelines-for-Assessment-and-](https://www.amstat.org/asa/education/Guidelines-for-Assessment-and-Instruction-in-Statistics-Education-Reports.aspx)
- 193 [Instruction-in-Statistics-Education-Reports.aspx\)](https://www.amstat.org/asa/education/Guidelines-for-Assessment-and-Instruction-in-Statistics-Education-Reports.aspx)
- 194 ● The Introduction to Data Science Curriculum ([www.introdatascience.org](http://www.introdatascience.org))
- 195 ● the draft *Data Literacy in K–12* (2020) and other resources from the Center for
- 196 Radical Innovation for Social Change (RISC) ([www.21cmath.org/](http://www.21cmath.org/))
- 197 ● The Messy Data Coalition
- 198 ● Youcubed data science resources, news articles, lessons and courses –
- 199 [www.youcubed.org](http://www.youcubed.org/datascience) /datascience
- 200 ● Statistical Literacy: A Complete Hierarchical Construct ([https://iase-](https://iase-web.org/documents/SERJ/SERJ2(2)_Watson_Callingham.pdf?1402525004)
- 201 [web.org/documents/SERJ/SERJ2\(2\)\\_Watson\\_Callingham.pdf?1402525004\)](https://iase-web.org/documents/SERJ/SERJ2(2)_Watson_Callingham.pdf?1402525004)

202 Two important sources for contexts in which to explore data science are

- 203 ● The *California Next Generation Science Standards* (CA NGSS) and
- 204 ● The California Environmental Principles and Concepts.

## 205 **The Statistical and Data Science Investigation Process**

206 The process of statistical and data science investigating is a four-part process:

### 207 **(1) Asking Questions**

208 Formulating questions that anticipate variability should be the beginning of the

209 investigative process. Examples of such questions include:

- 210 ● How fast will my plant grow?
- 211 ● Do plants exposed to more sunlight grow faster?
- 212 ● How does sunlight affect the growth of a plant?

213 These questions contrast with questions that are not investigative and have one

214 answer, such as: How tall is my plant? While questions start the investigative process,

215 students should be encouraged to ask questions throughout the investigative process.

216 (GAISE II page 15).



217 Recent work in the data feminism movement (see, for example, D'Ignazio & Klein,  
218 2020) has drawn attention to the need to understand not just the context of the data, but  
219 the motivation behind data collection and to ask questions about who has been included  
220 or excluded from data.

221 Survey questions will be important to students' investigations. These are questions  
222 designed to elicit data from people in order to address a statistical question, such as the  
223 length of time it takes to ride a bus to school.

224 As Arnold has stated, "Any question whose investigation requires repeated counting,  
225 measuring, or categorizing is one that data helps to answer." Students learn to use data  
226 in increasingly sophisticated ways. Early questions are primarily about description,  
227 beginning with categorizing and counting, expanding into questions in measurement  
228 situations (at first length/distance; later time, area, volume, and rates). Describing  
229 relationships between two varying quantities develops as students move through the  
230 grades, as do formal quantitative calculations.

231 CODAP provides a set of databases that will be interesting to school students, such as  
232 data on earthquakes, mammals, stars and cities, and an accessible data investigation  
233 online tool. Students can be encouraged to ask questions of the data. For example, a  
234 data set of mammals may raise the question, "Is the size of mammals related to the  
235 length of time they sleep?" Students can investigate questions using graphing tools that  
236 compare variables, statistical tools, a mapping tool and others (see  
237 <https://codap.concord.org/>).

238 Multivariate thinking comes naturally to humans, and students can develop curiosity  
239 about all sorts of data and situations. Young students may ask questions with one  
240 variable, such as what is the average age of my class? but as they get older we should  
241 encourage bivariate and multivariate thinking. Are older students at my school more  
242 likely to read more books would be an example of bivariate data collection.

243 Vignette 1: CODAP

244 A group of three students work to explore a CODAP database of 27 mammals:

245 <https://codap.concord.org/releases/latest/static/dg/en/cert/index.html?url=https://concord>  
246 [-consortium.github.io/codap-](https://concord-consortium.github.io/codap-)  
247 [data/SampleDocs/Science/Biology/27mammals/Mammals\\_Sample.codap](https://concord-consortium.github.io/codap-data/SampleDocs/Science/Biology/27mammals/Mammals_Sample.codap)

248 The database provides variables such as the height, mass, speed, life-span, and sleep

249 hours of the mammals. The students quickly become curious and ask questions like,

250 “Do bigger animals sleep longer?” They plot the two variables with the graph tool and

251 start to notice a relationship—in the opposite way than the one they thought—it seems

252 the bigger animals sleep less. The students start an animated conversation discussing

253 the reasons this might be, is it because they are more likely to be predators? They then

254 move on to investigate another relationship—who sleeps more, plant or animal eaters?

255 The students again notice a relationship as well as an outlier (the rabbit) so they wonder

256 about the rabbit, and look at more rabbit data. The students’ investigation of bivariate

257 data and their relationships is filled with moments of curiosity and excitement, as well as

258 important learning.

**Group 3: Tanya, Duane, Elexis, and Kevin**

**Our First Question: Do big animals sleep more than small animals? Does diet have anything to do with the size and amount of sleep an animal typically gets?**

**Our Second Question: Do plant-eaters tend to sleep less than meat-eaters or animals that eat both meat and plants?**

the rabbit is small in mass and sleeps an average amount compared to all other mammals

notice the yellow points!

We thought that big animals would need more sleep, but this graph shows us that some of the biggest animals sleep the least! We also noticed that it looks like plant eaters don't sleep much so we will explore that next

Yes! We can see that almost all of the plant-eaters represented in this data sleep less (on average) than the animals with other diets. There is one exception that sleeps more than the other plant-eaters!

when we compared the rabbit to other plant eating animals it stands out for sleeping much more than the others

There is one animal in the plant-eater group that is smaller than the rest, and it sleeps a lot more... we went back to the data and found that this animal is the rabbit!

259

260 **(2) Collecting and Considering Data**

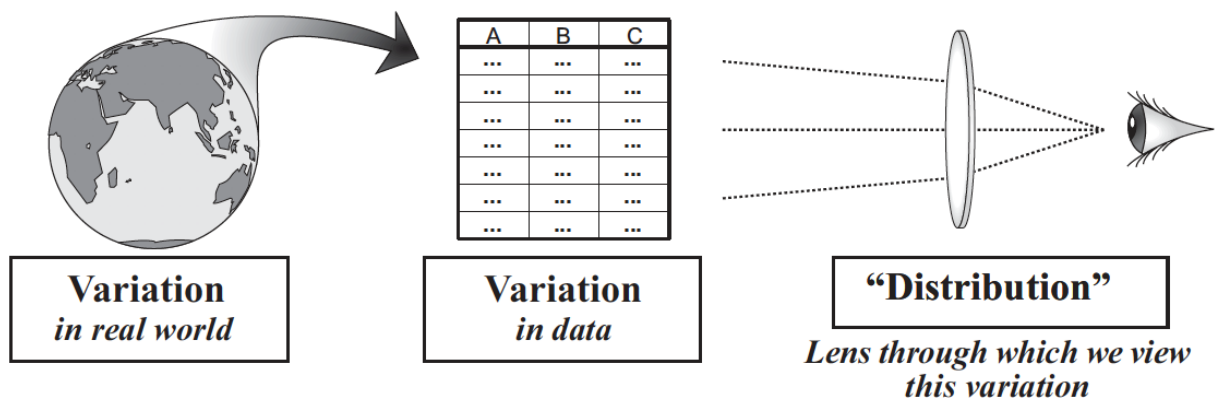
261 Sometimes students may collect their own data when investigating a question. For  
262 example, they may ask how far do students travel to school? or they may consider two  
263 variables, such as: Are students happier on sunny days? Or they may consider which  
264 plants are most prevalent in their local area. In all of these cases, students could collect  
265 data by observing plants or surveying students.

266 A key characteristic of data science is asking questions of “big data” – a data set that is  
267 complex, messy and includes many variables. Students can ask questions of big data  
268 sets and different students in a class may ask different questions. In a high-school data  
269 science class students can learn to clean data sets, an important part of the work of a  
270 data scientist. High school students can also learn to download and upload data—and  
271 develop the more sophisticated “data moves” that are important to learn if students are  
272 tackling real data sets.

273 After data is collected or acquired students should ask questions about the data—how  
274 do the variables differ? How were they collected? Who of what was included in the data  
275 collection. This helps students develop an understanding of variability.

276 High school students taking a course in data science may consider more complex  
277 conceptions of data science, that are located in the idea of variation, see for example  
278 Figure 2 from Wild (2006).

279



280

281 Further details of the understandings that may be developed in a high school course are  
282 given later in this chapter.

283 Sometimes students may be given data first and then ask a question of the data –  
284 reversing the order of 1) and 2).

### 285 **(3) Analyzing Data and Developing Meaning**

286 In the younger grades, students can analyze and develop meaning from data as they  
287 represent it in different ways, using picture graphs, line graphs, bar graphs and other  
288 forms of data visualization. From sixth grade, students can learn more formal methods  
289 to understand data. The field of statistics has been described as the study of variation,  
290 and students learn about variation when they receive opportunities to consider the  
291 distribution of data. Measures such as mean, median and mode are measures of the  
292 center of a distribution that students learn in middle school. CODAP tools allow students  
293 to see distributions of data and to see, visually, that the spread of a distribution will  
294 impact measures of center. In high school students will learn about measures of spread  
295 and about regression lines.

296 One of the features of data science is the possibility of predicting outcomes, such as the  
297 cable news programs' predictions of election outcomes. Developing understanding of  
298 what a prediction means, and how to compare predictive strength of one model over  
299 another is not simple and should be developed as a learning trajectory spanning several  
300 grades. Students who specialize in high school can learn about cross-validation  
301 techniques. Much of the work of professional data scientists is concerned with  
302 quantifying error from predictions.

### 303 **(4) Interpreting and Communicating Results**

304 Students learn to interpret data in increasingly sophisticated ways. Young students may  
305 make statements about their data or create data visualizations to communicate results.  
306 They may describe the difference between two groups. Even in the early grades,  
307 teachers can have conversations with students about generalizability—how much can  
308 we generalize from the data we have collected to broader populations? As students

309 move through the grades they can learn to generalize more formally and to include  
310 statements of probability and certainty.

311 A data scientist does not just perform calculation, and an important part of data science  
312 is the communication of results. Whereas statistics used to rely on bar charts, pie charts  
313 and other familiar representations, data science has created multiple forms of  
314 visualizations that represent data, as can be seen in Vignette 2.

315 Data science is about developing understanding of a situation, it involves holistic  
316 thinking, interpretation of meaning, and the communication of complex ideas. An  
317 effective data communication draws from writing, and visualizing as well as calculating.

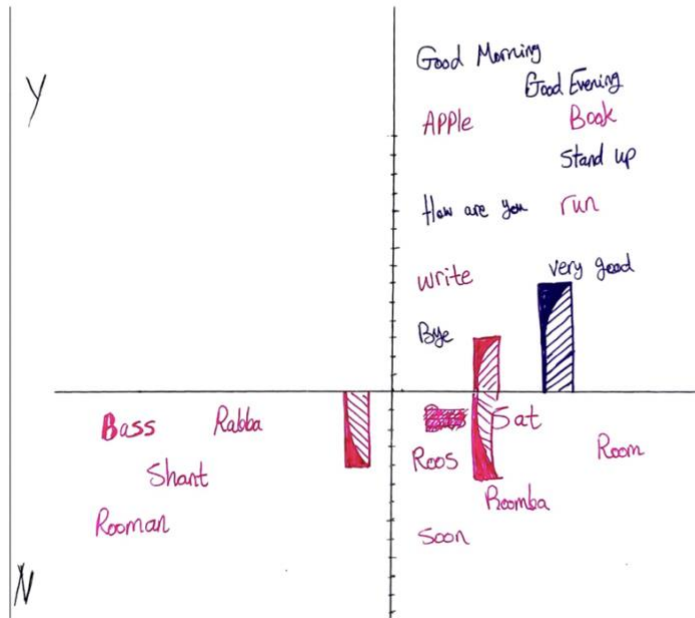
## 318 Vignette 2: Dear Data

319 Rico shares with his class of students the true story of two designers, who lived on  
320 different sides of the Atlantic Ocean—one in London, one in New York. For an entire  
321 year the two designers mailed each other a postcard every week, that included data  
322 from their lives, that they represented in creative and visual ways. The data  
323 representations included multiple variables. For example, some weeks the designers  
324 recorded all their moments of indecision, in another they recorded all the times that they  
325 laughed. The students looked at some of the data visualizations the designers produced  
326 and discussed what they could learn and how they could interpret the different variables  
327 (<http://www.dear-data.com/>)

328 After the discussion, Rico asked his students to collect data over at least a 24-hour  
329 period, collecting data on something that interested them, recording at least two  
330 variables. When the students came back to class with their data Rico organized the  
331 students into groups and asked them to create data visualizations together, supporting  
332 each other to consider ways they would represent different variables. In the discussion  
333 Rico payed attention to the language needs of the students, and the ways that the  
334 activity drew from the principles of Universal Design for Learning (UDL). Students were  
335 excited to make their data visualizations, such as the following:

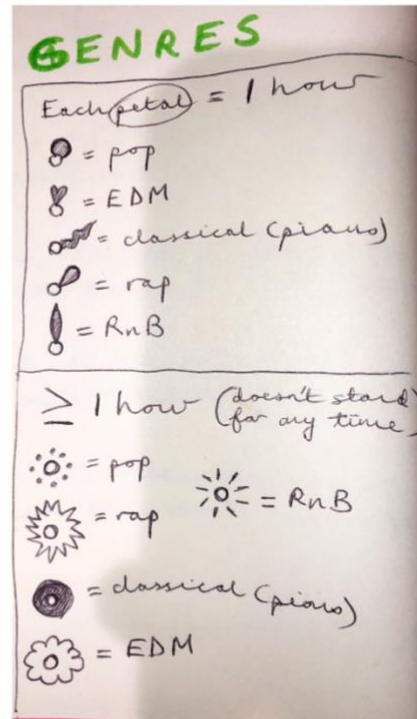
**Abdu**

How many times does Abdu's 6 year old sister use English and non-English words she knows or does not know while pretending to be a teacher?



**Nikita**

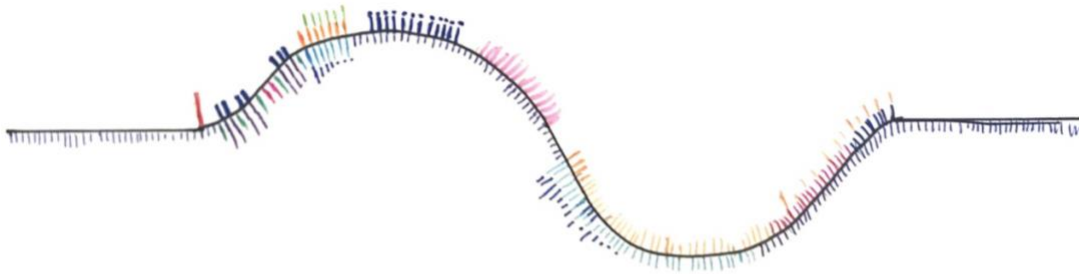
One week of listening to music, what type of genre it was and what Nikita was doing.





**Nathan**

Representation of sound length, level of loudness and how much attention is given to it.



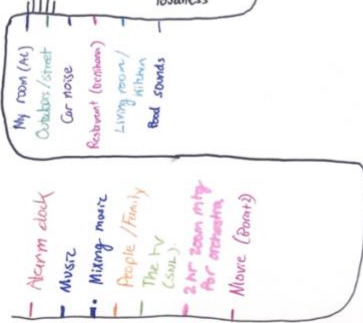
DEAR DATA

"SOUND WAVE"

Each line represents ~ 10 μm of sound.

Lines above are sounds you pay attention to  
Lines below are ambient

length of line = general loudness



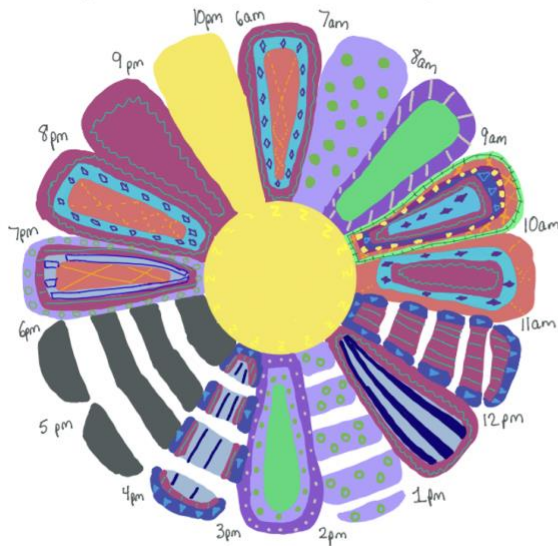


## Kira's dog interactions

Dear Data:

For a day, between 6am-10pm while I was awake, I recorded my interactions with my dog - Daisy, a golden doodle - and her (sometimes sassy) behaviors. Usually, she is my study buddy for the day.

This data is from Wednesday 11/4/20. Below you will find the key. Something to note, starting from the outermost layer of a petal and going inward accounts for the order of the actions.



My Actions	# of Occurrences	Daisy's Responses	# of Occurrences
<ul style="list-style-type: none"> <li>Physical Pet nr = belly rub    ■ = Paw</li> </ul>	7	<ul style="list-style-type: none"> <li>Roll over ↕ = from sitting    ↓ = from standing</li> </ul>	4
<ul style="list-style-type: none"> <li>Show a treat ☼ = cheese    ☼ = cookie</li> </ul>	2	<ul style="list-style-type: none"> <li>Walk away ■ = I annoyed her    □ = distraction</li> </ul>	4
<ul style="list-style-type: none"> <li>Call name ○ = actual name (course)    ● = nickname (Daisy, Doodle, etc.)</li> </ul>	4	<ul style="list-style-type: none"> <li>Come Solid fill = when called</li> </ul>	2
<ul style="list-style-type: none"> <li>Talk to Daisy x = scolding    x' = positive</li> </ul>	6	<ul style="list-style-type: none"> <li>Paw (beg) ▲ = unprompted    △ = prompted</li> </ul>	3
<ul style="list-style-type: none"> <li>No interaction</li> </ul>	2	<ul style="list-style-type: none"> <li>Sleep Z = herbal    Z' = other</li> </ul>	16
<ul style="list-style-type: none"> <li>Gave Daisy a shower □ = Paws    ■ = full</li> </ul>	1	<ul style="list-style-type: none"> <li>Muddy # = digging    = = rain</li> </ul>	1

Not in class = solid    While in class = dashed

339

340 The students made their visualizations using Google Jamboards. After they made them  
341 Rico asked the groups to look at the work of other groups and provide feedback to each  
342 other on a sticky note. The students were excited to see the ways the different variables  
343 related to each other and the ways they could be represented.

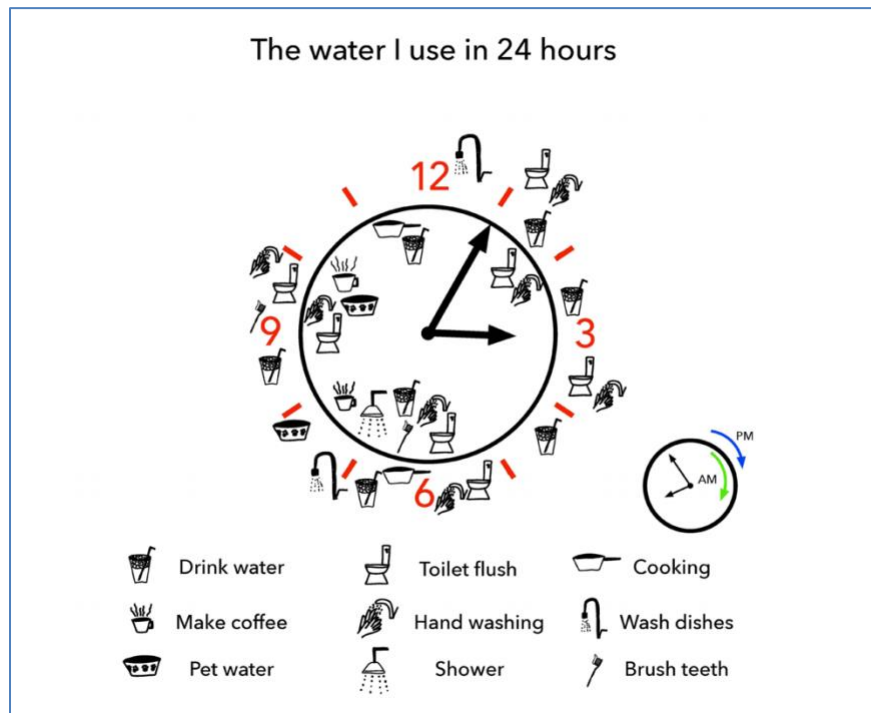
### 344 Data Talks K-12

345 Data talks are short classroom discussions to help students develop data literacy. This  
346 pedagogical strategy is similar in structure to a number talk, but instead of numbers  
347 students are shown a data visual and asked what interests them. The idea of a data talk  
348 was inspired by a New York *Times* weekly section called, "What's Going on in this  
349 Graph?" In the New York Times, students can submit their own ideas to a member of  
350 the American Statistical Association, who reveals their thinking on the data in the  
351 graphs. In the classroom the teacher can guide the discussions and help students

352 develop important understandings. However, it is important to recognize that teachers  
353 do not have to be an expert in the topic of the data visualization—instead teachers can  
354 guide and encourage curiosity and question asking. One way to support thinking and  
355 speaking like a mathematician is to incorporate writing activities or math journals, which  
356 allow students to process learning and continue questioning. These activities help all  
357 students gain and exchange information and ideas, and support the *California English*  
358 *Language Development Standards*’ three communicative modes (collaborative,  
359 interpretive, and productive), and allow them to apply knowledge of language to  
360 academic tasks using various linguistic resources.

361 If questions cannot be answered by the teacher or students they can be investigated  
362 further. Data talks are intended to pique students’ curiosity and encourage question  
363 asking, and to help them understand and “read” the data-filled world in which they live.  
364 Many of the data visualizations illustrate how multiple variables can be incorporated into  
365 one graphic—allowing students to think multivariately.

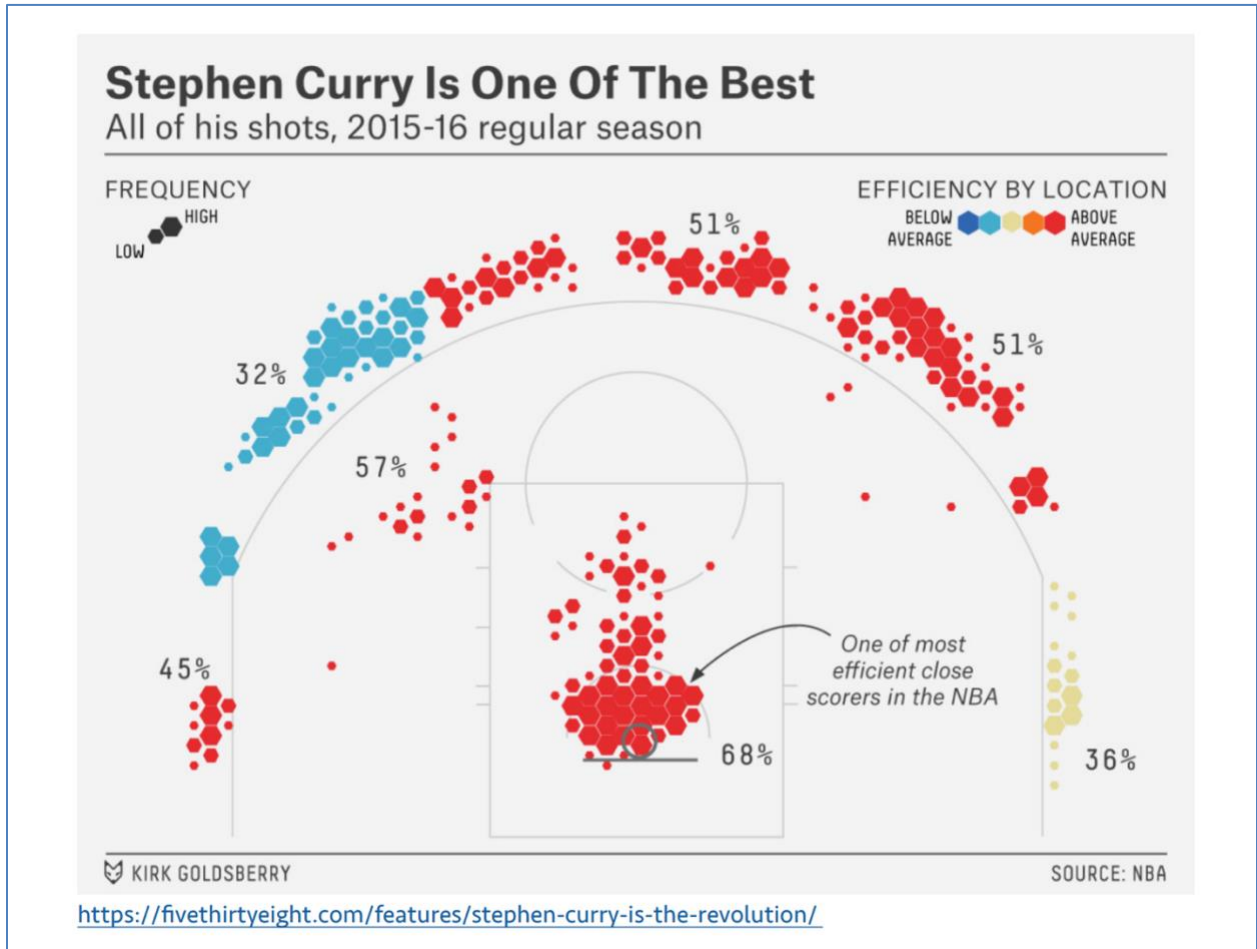
366 Grades with younger students can use data visualizations with no or few numbers, or  
367 smaller numbers for example:



368

369 Source: <https://www.youcubed.org/wp-content/uploads/2020/09/Water-Usage.pdf>

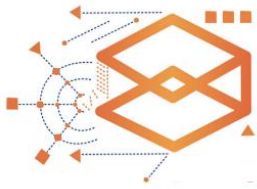
370 From grade five, students should be able to interpret data visualizations with  
371 percentages, for example:



372

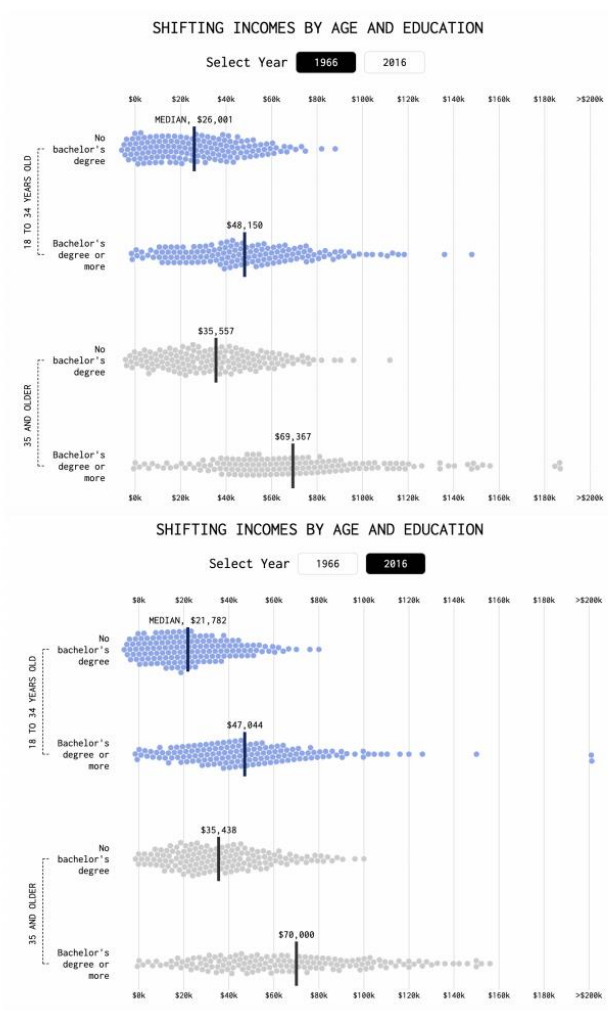
373 In higher grades data visualizations can include more complex data representations.

374 For example:



## Youcubed Data Talk Shifting Incomes

What do you notice?  
What do you wonder?  
What is going on in this data visualization?



<https://flowingdata.com/2017/05/02/shifting-incomes-for-young-people/>

375

376 The above examples, and more, can be found at Youcubed,

377 (<https://www.youcubed.org/resource/data-talks/>), Flowing Data

378 (<https://flowingdata.com/>), Five Thirty Eight (<https://fivethirtyeight.com/>), and the New

379 York *Times* resource itself ([https://www.nytimes.com/column/whats-going-on-in-this-](https://www.nytimes.com/column/whats-going-on-in-this-graph)

380 [graph](https://www.nytimes.com/column/whats-going-on-in-this-graph)).

## 381 Transitioning from Pre–K

382 Before kindergarten, children begin to describe their world in language, identifying  
383 characteristics of objects, places, people, and events: *The ball is red. My classroom is*  
384 *warm. My teacher is in their twenties. Our trip to the park was too short.* Identifying  
385 characteristics is the beginning of data, and wondering about characteristics—including  
386 countable characteristics—is the beginning of asking questions that data can help to  
387 answer. In the California Preschool Learning Foundations, this content is located under  
388 the heading of “Algebra and Functions (Classification and Patterning),” in which children  
389 “sort and classify objects in their everyday environment,” (by one attribute at around 48  
390 months and by more than one attribute at around 60 months of age); and in  
391 “Measurement,” in which students compare and order objects directly at around 48  
392 months of age and may use an intermediate object to compare at around 60 months of  
393 age (Preschool Learning Foundations, Volume 1). These preschool activities directly  
394 enable the types of kindergarten through grade five learning trajectory described below.

## 395 K–5

396 The big ideas of data in these early grades include

- 397 ● Data for understanding. What questions can we ask? What data do we need to  
398 answer it?
- 399 ● Defining data: What is data and how where data collected?
- 400 ● Representing and interpreting data: What does data look like and what does it  
401 mean?

402 These ideas are represented in the most important pedagogical and practical process  
403 through which data plays a role in making sense of the world, outlined in these four  
404 steps, from GAISE II and discussed above.

- 405 1. Ask a question (SMP.1: Make sense of problems and persevere in solving them)
- 406 2. Collect and consider data
- 407 3. Analyze data and develop meaning (SMP.2, SMP.4, SMP.7)
- 408 4. Interpret and communicate results (SMP.4, SMP.5)

409 An important distinction to consider is between *categorical* (non-numerical) data and  
410 *measurement* or *quantitative* data. For instance, consider a set of colored blocks in the  
411 classroom. “Color” is a categorical variable students could observe about each block.  
412 “This block is 15 centimeters long” is a measurement data point. The standards develop  
413 categorical data in grades K–3 and measurement data beginning in grade two. A  
414 description of the development of the data and measurement standards organized  
415 according to *categorical* vs *measurement* data is found in the *Measurement and Data*,  
416 *K–5* progression document. [Note: link to the Progressions documents currently is:  
417 <http://ime.math.arizona.edu/progressions/>; a new link will be provided on the CDE’s  
418 website in future drafts]

<p><b>Categorical data</b></p> <ul style="list-style-type: none"><li>• Color (red, green, blue, yellow) of blocks in the class set</li><li>• Species of trees on the school grounds</li><li>• Identification of schools in the district as “elementary school,” “middle school,” or “high school.”</li></ul>
<p><b>Quantitative (or Measurement) data</b></p> <ul style="list-style-type: none"><li>• Height (or circumference of trunk, or biomass) of trees on the school grounds</li><li>• Number of pages (or weight, or height) of books in the classroom</li><li>• Annual income for households in a census tract</li></ul>

419 Figure 1: Examples of Categorical and Quantitative Data

420 What questions can data help to answer?

421 All work with data should begin with noticing and wondering: “I notice that...” or “I  
422 wonder what...” or “I wonder how many....” (see <http://www.nctm.org/mathforum/>) To  
423 prompt wonder, teachers can ask: “What do you notice or wonder about here [in this  
424 context], that we could (count/measure/keep track of) to figure out or explore further?”  
425 To establish effective routines, and to support language development in “I wonder”  
426 activities, it can be effective to provide these examples as sentence starters.

427 As students gain confidence in their ability to speak like mathematicians, statisticians  
428 and/or data scientists, the teacher should encourage students to generate questions

429 themselves to build their agency in using mathematics to make sense of their worlds. A  
430 weekly whole-class “I wonder” routine—in which students propose questions to  
431 investigate by collecting data—would build a powerful practice of observing the world  
432 with a data lens, contributing to students’ development of modeling with mathematics  
433 (SMP.4).

434 In kindergarten, students compare the number of objects in different categories  
435 (K.CC.C.6) to answer “Which has more?” questions (*I wonder whether there are more*  
436 *square blocks or more triangular blocks on the desk?*). At first, the teacher suggests or  
437 specifies categories; eventually students generate ideas for classification. They also  
438 directly compare (as opposed to measuring with a unit or an intermediate) objects with  
439 common measurable/countable attributes to see which has more (K.MD.A.2, K.G.B.4) (I  
440 wonder which shape has more sides?; Which kind of block is heaviest?—using a  
441 balance or informal one-in-each-hand comparison, rather than a scale). “I wonder...”  
442 questions should explore both of these: two-category, “Which is more?” questions and  
443 comparison of objects according to length, height, weight, and countable attributes like  
444 number of sides. Student-generated questions provide opportunities to work on  
445 precision of language as well—for example, when students are asked to clarify what  
446 they mean by “bigger.” Mathematics discussions that are rooted in academic language  
447 will help students understand mathematical concepts more deeply as well as discover  
448 new ones. As the years progress, students or teachers may reach beyond the  
449 classroom to find contexts for their: “I wonder...” questions.

450 In addition to questions that can be answered with a single value, students can start to  
451 pose statistical investigative questions that involve multiple variables such as, I wonder  
452 if plants grow more with more sunlight? Or I wonder if age affects which color people  
453 like?

454 In first grade, measurement of length and time are the contexts to emphasize in  
455 generating questions (1.MD.A.2, 1.MD.B.3), along with continued work categorizing and  
456 counting objects (1.MD.C.4) and categorizing geometric objects by attributes (1.G.A.1).  
457 Second graders should continue to explore questions in length measurement  
458 (2.MD.D.9) and time (2.MD.C.7) contexts, and add money contexts (2.MD.C.8). When

459 selecting “I wonder” questions, it is important to avoid situations that serve as markers  
460 for economic or social status, e.g. “I wonder who has the most expensive backpack,” “I  
461 wonder who is the most popular kid in school,” or, perhaps less obvious, “I wonder who  
462 has the newest shoes.” It is similarly important to avoid questions about students’  
463 physical attributes, even those that seem innocuous such as height or arm length.  
464 Instead, some good questions to wonder about might be “I wonder what time it will be  
465 when the next person walks into the classroom” or “I wonder which book in the  
466 classroom is the most read,” comparing events or objects rather than personal  
467 characteristics.

468 In third grade, contexts for questions to investigate using data should expand to include  
469 volume and mass measurement (grams, kilograms, and liters, but not compound units  
470 such as  $\text{cm}^3$ ) in addition to the length, time, and money contexts from earlier grades  
471 (3.MD.A.2). Time measurements are refined to the nearest minute (3.MD.A.1) and  
472 length now includes half- and quarter-inches (3.MD.B.4). Beginning ideas of area give  
473 another possible context, limited here to areas that can be covered by a whole number  
474 of unit squares (3.MD.C.5, 3.MD.C.6).

475 In fourth grade, a significant context for data-investigation questions is classification and  
476 analysis of two-dimensional shapes (4.G.A.2). Incorporating this Geometry standard to  
477 help build data understanding can foster the important practice of analyzing by  
478 attributes—one instance of SMP.7 (Look for and make use of structure). Fourth-grade  
479 students also extend the set of units they work with (4.MD.A.1) and can generate data  
480 about area for more complex shapes. Fifth graders deepen their understanding of  
481 volume to include unit cubes, making this an important context for data-inquiry  
482 questions. A teacher could invite students to build a structure out of multi-link cubes and  
483 then collect data from the class by asking, for example, how many cubes they use in  
484 each of their different structures they built, or the height and width of their structures,  
485 and color of the blocks. Students can collect data on multiple variables.

486 In K–5, “I wonder...” questions come primarily from personal experience. See below for  
487 additional examples.



## 488 Asking Questions, Collecting and Analyzing data

489 Questions invite inquiry. An important part of students' K–5 experience should involve  
490 coming to recognize that, when they choose and pose questions, they can collect or  
491 analyze data to find answers (SMP.4). Some of the most valuable conversations about  
492 data occur when students notice patterns in a data set and begin asking questions.  
493 Remaining alert for these everyday moments—perhaps in attendance, weather, or  
494 lunch-count data—may generate opportunities for discussing statistical investigative  
495 questions and exploring how data can help answer them.

496 As students come to pose authentic questions such as the ones described above, they  
497 should also encounter opportunities to help determine how data might be produced to  
498 answer them. In addition to producing data directly through their own observations,  
499 students should gain exposure to designing and using surveys and simple experiments  
500 to generate data. By producing their own data from their classroom or community (*How*  
501 *does age of students relate to their enjoyment of school? Does time on social media*  
502 *apps increase with age? How much waste is generated by different companies/our*  
503 *school?*), students recognize data as having context and deriving from observation and  
504 measurement, and they come to see data (and mathematics more broadly) as a tool to  
505 help think about their worlds. Data gathered by others (such as those in the data talks)  
506 can help to answer questions students generate about their own communities.

507 When choosing data tasks that include categorizing and counting, consider the grade  
508 level expectations for counting (up to 10 objects scattered, or up to 20 if arranged in a  
509 line, array, or circle, in kindergarten [K.CC.B.5], 120 by the end of first grade  
510 [1.NBT.A.1], and up to 1000 by the end of second grade [2.NBT.A.2]). Such tasks can  
511 also be structured to build place value understanding.

512 In kindergarten, once students notice things in a context and wonder about a question,  
513 they describe measurable, countable, and observable attributes of objects or situations  
514 (K.MD.A.1, K.G.A.1, K.G.B.4), and classify objects and count the number in each  
515 category (K.MD.B.3), such as categorizing a set of cubes by color. In this last context,  
516 both “this cube is red” and “there are 13 red cubes in the set” are data points. Notably,

517 most work on *number* in kindergarten should be with numbers representing quantities of  
518 objects (SMP.2); thus, most numbers encountered in kindergarten are actually data.

519 In first grade, students explore their time and length questions by measuring lengths of  
520 objects which are a whole number of units (1.MD.A.2) and telling and writing time in  
521 hours and half-hours. Counting and categorization situations should include up to three  
522 categories (1.MD.C.4). Second graders measure length to the nearest whole unit  
523 (2.MD.D.9), using different standard units (centimeters, meters, inches, feet) (2.MD.A.3)  
524 and several tools (2.MD.A.1) and measure time to the nearest five minutes (2.MD.C.7).

525 Students in third through fifth grades refine their measurements of lengths and time, and  
526 expand the set of units they use; and they add area and volume measurement to their  
527 repertoires (as described above in “What questions”). By the fifth grade, students should  
528 understand that data sets can include different types of variables, such as categorical  
529 and quantitative. They should recognize that an individual instance or object can  
530 possess attributes that exemplify these different types, and should have gained  
531 experience measuring, characterizing and analyzing such diverse types of data and  
532 associating them together.

533 An important understanding that students need to develop through grades K–5 is the  
534 idea of variability and variables. When students ask questions such as: How high are  
535 the plants in the classroom? they are considering one variable: height. When they  
536 consider whether older students spend more time on social media apps they are  
537 collecting bivariate data—with two variables—age and time. When they make their own  
538 data visualizations, as seen in Vignette 2, they may collect data on multiple variables.  
539 Multivariable thinking is important to develop through the grades.

## 540 Interpreting and Communicating Results

541 Sorting objects into two categories and representing these categories by their count  
542 (K.MD.B.3) is a first example of students representing data to help make sense of their  
543 worlds (SMP.4). First grade students organize up to three categories and ask and

544 answer questions about the relative sizes of categories and about the total number of  
545 data points.

546 Second grade begins an expanded focus on data representation, introducing line plots  
547 (whole number units only; 2.MD.D.9), picture graphs, and bar graphs. These graphs can  
548 be used to answer put-together, take-apart, and compare questions (2.MD.D.10). In  
549 third grade, *scaled* picture and bar graphs are added as a tool for visualizing “how many  
550 more” questions (3.MD.B.3), and line plots may have half-unit and quarter-unit markings  
551 as appropriate (3.MD.B.4). In fourth grade, line plots may display additional fractional  
552 units (to eighth-units), and be used to answer additional questions about differences—  
553 between maximum and minimum measurement, for example.

554 Fifth grade does not extend the expected set of data representations, but students do  
555 use line plots in a sophisticated way that sets the stage for understanding the most  
556 common measure of *center* for a set of data—the *mean* (commonly called the  
557 average)—in sixth grade. Namely, fifth grade students use a line plot to decide how a  
558 repeatedly-measured quantity could be redistributed equally (5.MD.2): “Given different  
559 measurements of liquid in identical beakers, find the amount of liquid each beaker  
560 would contain if the total amount in all the beakers were redistributed equally.”

561 While the data visualizations mastered by fifth grade only include picture graphs, bar  
562 graphs, and line plots, students do not need to be restricted to these. Each of these  
563 represents repeated measurements of a *single* varying quantity; science curricula in  
564 particular, and many questions of interest in general, require the consideration of  
565 relationships between *two or more different* changing quantities, such as erosion and  
566 time (NGSS 4-ESS2-1 Earth’s Systems) or length or direction of shadows and time  
567 (NGSS 5-ESS1-2 Earth’s Place in the Universe). Such reasoning involving multiple  
568 variables is an important aspect of modern encounters with data, and students should  
569 experience it at all levels. Although the scatter plot, a crucial data representation tool for  
570 two varying quantities, is not mastered until eighth grade (8.SP.1), it must be explored  
571 informally much earlier for students to be able to meet the eighth-grade expectations.  
572 For example, students can plot quantities changing over time (e.g. height of a plant,  
573 length of the day, high temperature for the day, temperature of a glass of water every

574 minute for an hour), with time on the horizontal axis and the changing quantity on the  
575 vertical. Once such a plot is created, it is an excellent context for a “notice and wonder”  
576 discussion.

577 In recent years, new technological tools and developments in data science have  
578 prompted an explosion in interesting data visualizations, many of which are quite  
579 comprehensible to young students with some exploration. Experiences with different  
580 visualizations will further expand students’ sense-making opportunities and encourage  
581 them to think about what they can understand looking at data sets in different ways. The  
582 examples from the New York *Times*, Youcubed, and other places illustrate multivariate  
583 data displayed in creative ways. This example shares the most popular songs each  
584 summer: <https://www.youcubed.org/resources/whats-going-on-in-this-graph/>.  
585 Newspapers and online news sources offer other examples; student-gathered examples  
586 help to build buy-in for a “can we figure out what this visualization is trying to help us to  
587 understand?” routine.

588 Interpreting data is a matter of making inferences from the data available. While  
589 students will encounter quantitative and nuanced techniques for making inferences in  
590 later grades, they should nevertheless encounter opportunities to make claims and infer  
591 conclusions across their K–5 years (SMP.3). When they do, students should learn both  
592 to wonder whether patterns or trends they notice in data extend beyond the particular  
593 group that generated the data, and to be skeptical about such extensions to larger  
594 populations (including considering ways in which the group might not be representative  
595 of the larger population). Additionally, students should learn that good claims draw upon  
596 data as evidence and that they always come hand in hand with a degree of uncertainty.  
597 Modeling the use of appropriate terminology such as “tends to,” “typical,” “usually,” and  
598 “similar” can help lay important groundwork for this concept (Rugin, 2019).

599 Preparing for the major data science work of grades 6–8

600 **Understanding Variability:** Variability is everywhere, and understanding variability is  
601 the core of developing data sense. While understanding of statistical variability and  
602 distributions is not in the standards until middle school, it is essential that K–5 students

603 encounter many experiences with variation, including counting, measuring, and  
604 observing quantities and characteristics that vary in order to be prepared for the first big  
605 idea in the Grades 6–8 section below. In particular, their encounters with data  
606 representations should highlight important ideas that set the stage for more involved  
607 work with distributions. When working with visualizations of data, students should  
608 consider not only the most popular value in a dataset (the mode) but also describe the  
609 shape and spread of data distributions. Identifying the maximum and minimum values of  
610 quantitative datasets can help students appreciate the concept of range as a measure,  
611 and looking for clusters and gaps in a distribution can begin to help them attend to its  
612 shape. As they engage in experiences where they produce their own data through  
613 measurement, teachers should highlight for students the variation that results.  
614 Measuring the same variable on multiple individuals or objects, for example, results in  
615 data that vary, and students should consider the causes or sources that might have  
616 given rise to the variation they have observed, working as they do so to differentiate  
617 between variation and error. For example, if students plant a particular variety of flower  
618 seed at multiple locations around the school, then measure the plants' height and the  
619 amount of sunlight each month, they can conduct investigations into the ways plant  
620 growth and sunlight relate to each other. They should discuss and describe any patterns  
621 in their bivariate data, and discuss reasons for the variability. Finally, they should  
622 consider their own measurement techniques, and how confident they are that they all  
623 measured the same way (so that if someone else measured, they would get the same  
624 height or sunlight).

625 **Randomness, probability, and uncertainty:** Randomness is a complex idea  
626 encompassing uncertainty *and* a level of predictability. When (blindly) drawing a cube  
627 out of a bag containing three blue cubes, two red cubes, and one yellow cube, nobody  
628 can predict with certainty what will happen on a single draw. But, over many draws, the  
629 person who always predicts a blue cube will be right about half the time. Activities that  
630 demonstrate this can be used to generate data for many of the explorations of the big  
631 ideas above, which will leave students well-prepared for a more formal treatment of  
632 randomness and probability in middle school. At this point, students should begin to

633 conceive of probability as a measure of the chance that something will happen, seeing it  
634 as a basic measure of certainty or uncertainty.

635 **Technology:** California’s 2018 K–12 *Computer Science Standards* include computer-  
636 based data sorting, categorizing, and visualizing for students in grades K–2 and 3–5  
637 (CS K–2.DA.8, CS K–2.DA.9, CS 3–5.DA.8). These standards are important  
638 preparation for middle and high school use of data software to visualize and interpret  
639 large data sets.

640 Finally, it is worth noting that (as in science and other fields) many questions that  
641 students might wonder about will not be fully answerable using K–5 tools. It is important  
642 that teachers have resources for helping students figure out which aspects of questions  
643 can be investigated with currently available tools, and have some understanding of data  
644 science tools which students will encounter later. For example, many will wonder about  
645 relationships between two different variables: *If I get up earlier, do I feel tired earlier in*  
646 *the afternoon at school? Do students who skip lunch eat more candy in the afternoon?*  
647 When one of the variables is categorical (like the skipping lunch question), separate line  
648 plots can be made for each category and the line plots compared. When both variables  
649 are quantitative, students could input data into CODAP and investigate the relationships  
650 by plotting their data on graphs, observing their distributions, and adding line plots.  
651 Another option is that one of the variables can be made into a categorical variable by  
652 defining categories in terms of the quantitative variable. For instance, waking-up times  
653 could be classified into “early” and “late” (ideally with a student-generated cut-point  
654 between early and late) and then dot plots of “time in the evening when I felt tired”  
655 created for each category.

## 656 Vignette: Logan from Kindergarten through Grade 5

657 A small sampling of Logan’s data science experiences in grades K–5 is described  
658 below. This is not intended to capture *all* of their data science experience, only to  
659 indicate a development towards powerful uses of data to understand their world. In each  
660 grade, Logan generated questions and gathered data (steps 1–3 in the process

661 described at the beginning of this K–5 section) and represented and interpreted data  
662 (steps 4–5).

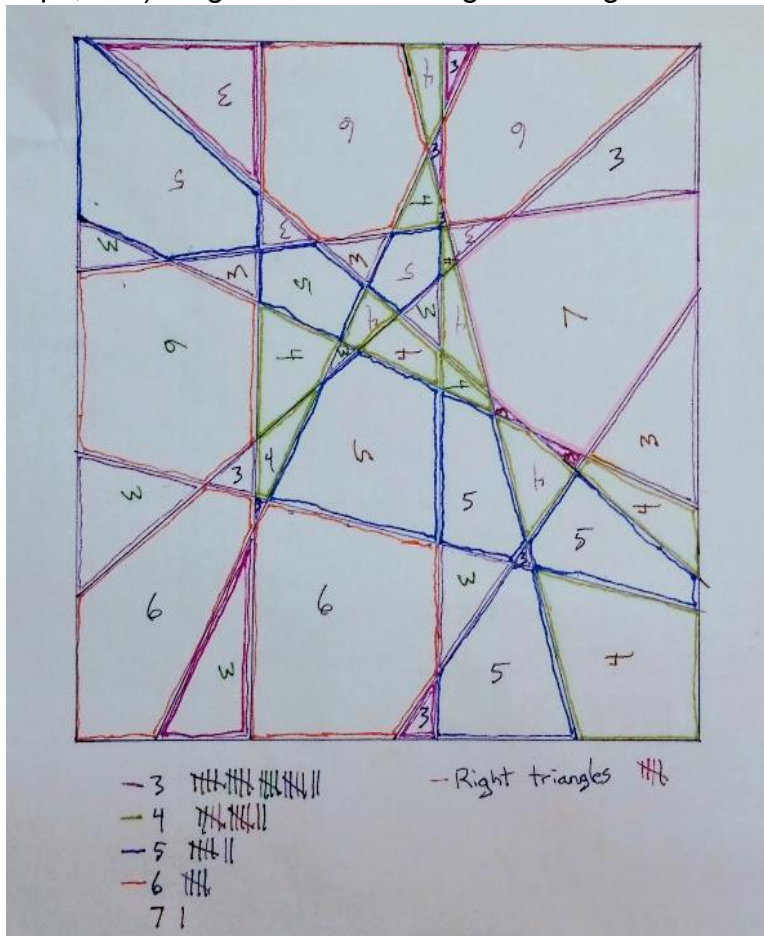
663 Logan entered kindergarten as a very active child, and always gravitated towards the  
664 group of children who ran the longest and climbed everything. Logan’s teacher asked  
665 lots of questions of students about what they noticed in the classroom and around  
666 school, both inside and out. These ranged from specific “how many in each category”  
667 questions (how many classroom doors of each color are there?), to direct comparison  
668 questions (which slide is taller, the blue one or the green one? How do you know?), to,  
669 eventually, *types* of questions: What are some things at school whose size we could  
670 compare? Recording students’ observations and category counts allowed all students to  
671 pose and answer “relative size” questions.

672 In first grade, student teams were asked to think of two similar things at school, such  
673 that they weren’t sure which was taller, and then to find a way to compare their heights.  
674 A variety of materials was available to use in the comparison. Logan’s team was able to  
675 compare the height of the slide in front of the school with the height of the slide behind  
676 school, measuring the height of both using towers of large DUPLO® bricks. The whole  
677 class used their data to discuss how much taller the slide in front was. After this, Logan  
678 wanted to build DUPLO® towers to measure height and length of lots of things, and was  
679 disappointed that the class didn’t have enough bricks to measure the height of the  
680 school (and that their teacher wouldn’t let them climb the school). As a class, students  
681 checked the length of the day (sunrise to sunset) each day, and maintained a running  
682 visible tally of the number of school days with less than 11 hours of daylight, 11 to 13  
683 hours, and more than 13 hours, for the entire year. As a class, they discussed what they  
684 thought might happen to the number of hours of daylight in the future and checked the  
685 data a month later to see whether their predictions were correct.

686 Logan’s second grade class marked their own yard sticks (marking a wooden blank in  
687 inches using only a 3-inch by 5-inch card), and then used them extensively to measure  
688 objects of interest to the nearest inch. Later they added centimeter markings to the  
689 other side of the yardstick, and discovered that measuring the same things with smaller  
690 units led to larger number measurements. When choosing an activity to time, Logan’s

691 group decided to time and record the amount of time in a week that team members  
692 spent reading in school, and compare those measurements over several weeks (this  
693 had the benefit that team members read much more during those weeks!). Other teams  
694 measured time spent playing outside, listening to announcements, and working at math  
695 stations. Teams made line plots of their data, and compared the line plots of different  
696 activities to discuss how students typically spend their school time.

697 As mass and volume became available in third grade as characteristics to measure,  
698 Logan's class used length/height, mass, and volume measurements to examine  
699 collections of objects. The line plots of the masses and the lengths/heights of the  
700 objects in the science corner looked quite different from each other; similarly for the line  
701 plots of volume, height, and mass of all objects in the room which hold water (vases,  
702 cups, etc.). Logan's team had a great disagreement about whether a taller vase should



703 hold more water than a shorter vase; the class eventually decided that this was usually  
704 but not always true.



705 One of Logan’s favorite activities in fourth grade was one that combined data work with  
706 classifying shapes by attributes: collaborative art pieces: Each team had a  $\frac{1}{2}$  meter by  
707  $\frac{1}{2}$  meter square on the board, and each student in the team drew in two edge-to-edge  
708 straight lines of their choice, using their meter sticks. Then one student in class chose a  
709 shape to try to find in the drawings, and each team outlined each new instance of that  
710 shape they found and described how they knew it was a (triangle, rectangle, right  
711 triangle, quadrilateral, etc.); this was repeated for several other shapes. They made an  
712 individual card to represent each piece of artwork, using the card to represent the  
713 artwork the different variables they measured for each piece (how many triangles, how  
714 many instances of each color, how clean or messy each line was, etc.) When they had  
715 made a full set of cards, they sorted them in various ways, then made a table to  
716 compare the tallies for the different pieces, discussing how the different features of the  
717 art and the process of creating it that might help explain the variations in their data.

718 By fifth grade, Logan and their classmates had constructed many line plots, and thus  
719 often wondered about quantities that vary on repeated measurement: The cartons of  
720 milk from lunch say they each contain 8 fluid ounces, but yours feels heavier than mine;  
721 am I getting less or are you getting extra? The weather site says the average high  
722 temperature here is 57°F (degrees Fahrenheit) in November, but today it got up to 65°F.  
723 How can we check whether this month is near average? To explore the first, the school  
724 donated 20 cartons of milk to the experiment When they examined the dot plot of their  
725 measured volumes, they saw that it had a tightly clustered shape, with a minimum  
726 measurement of 7.8 fluid ounces and a maximum of 8.2 fluid ounces, and that the most  
727 frequent value was 7.9 ounces. One student in the group thought that some milk  
728 probably remained in the containers, so the group spent a while trying to figure out how  
729 they might identify how much had been left inside (teams came up with several  
730 methods). For the second, the class recorded the daily high all month, recorded them  
731 on a line plot which also had marked the “average” high temperature from the weather  
732 site, and used the line plot at the end of the month to discuss whether it was consistent  
733 with the stated average (without computing an average of the data).

734 Students like Logan, with a rich variety of experiences using data to explore contexts  
735 and questions of interest, will be well-prepared to use mathematics, and data in  
736 particular, to make sense of their ever-expanding worlds. Data sets will get larger, and  
737 contexts wider, in ensuing years.

## 738 Grades 6–8

739 Middle school includes a big expansion in important ideas. The big ideas of data  
740 science include

- 741 • Data in the world: exploration, interpretation, decision making, ethics
- 742 • Statistical variability: Describing, displaying, and comparing
- 743 • Sampling to understand a population: randomness, bias, how many?
- 744 • Are they related? Multivariate thinking
- 745 • What are the chances? Probability as the basis for data-based claims

746 As in earlier grades, students experience data science as a tool to help understand their  
747 worlds via a process that begins with wondering questions. This is also the beginning of  
748 the mathematical modeling cycle (Pelesko, 2015) and the statistical and data science  
749 exploration process, and of investigations in science (NGSS Lead States, 2013).

750 The GAISE II Statistical and Data Science Exploration Process:

- 751 1. Curiosity and question asking
- 752 2. Collect and consider data
- 753 3. Analyze data and develop meaning
- 754 4. Interpret and communicate results

755 This process, beginning with noticing and wondering, often gets lost in the details of  
756 step 3, which contains the different statistical methods that have been developed for the  
757 analysis of data. It is crucial to keep all work with data tied to authentic questions.

758 Prediction is a key activity that builds student ownership of the process and conclusions;  
759 it also builds a habit of asking “does this make sense?” (SMP.1) by comparing results  
760 with expectations.

761 Data in the world: Question asking, exploration, interpretation, decision  
762 making, ethics, technology

763 What functions does data science play in the modern world?

- 764 ● **Question Asking and Exploration:** Data science and statistical exploration  
765 starts with questions that are posed by students. When students are invited to  
766 wonder about situations, and when they are given interesting datasets they will  
767 become curious and can ask questions of data that they can explore and  
768 investigate. Data exploration includes understanding the context and situation,  
769 data should never be abstracted from their context. Students can look for hidden  
770 patterns and associations. Any patterns or associations discovered can lead to  
771 new conjectures or questions to investigate further. In eighth grade, students can  
772 begin this process with datasets that include multiple variables, such as those  
773 given in CODAP (<https://concord-consortium.github.io/codap-data/>). As  
774 mentioned in the chapter introduction, vast quantities of data are collected every  
775 day, and only a small fraction are analyzed.
- 776 ● **Interpretation:** Every encounter with data should revisit the context from which  
777 the data originated, interpreting results of data analysis in that context. This  
778 includes answering any questions that began the encounter and reporting any  
779 other associations or patterns that were discovered.
- 780 ● **Decision making:** Commonly, data is used to inform decisions following the  
781 question/data/represent/interpret process. Often, however, data is used to justify  
782 and explain a decision, even if data didn't play a meaningful role in the decision.  
783 There is great potential for abuse here, by including data collections that support  
784 the predetermined decision and leaving out those that do not.
- 785 ● **Ethics:** Modern ubiquitous data collection raises a host of ethical questions, both  
786 about how and what data is gathered and stored, who is included or excluded in  
787 the data, and how that data is used and presented. Middle school students need  
788 to understand their own online data footprint (for example, how companies  
789 aggregate information about individuals to create detailed profiles) and should

790 confront scenarios in which they must make decisions in hypothetical situations  
791 involving data exposure, consent for data collection, etc.

792 ● **Technology:** California’s 2018 K–12 *Computer Science Standards* (CSS) expect  
793 students to make use of computers for data organization and visualization (CSS  
794 6–8.DA.8). More importantly, given the amount of data collected and stored  
795 today, real-world datasets are incomprehensible without such computer  
796 assistance. Students should use modern data software extensively, especially for  
797 organizing and displaying features of data set.

## 798 Describing, displaying, and comparing statistical variability (Grades 6–7)

799 Sixth-grade students build on earlier experiences by distinguishing between statistical  
800 questions, which can be investigated using data that varies (analysis of social media  
801 usage by age of students), and questions without variations in (correct) responses (How  
802 many days are there in January?) (6.SP.A.1). When considering a statistical question,  
803 they understand that the variation in numerical data has a distribution which can be  
804 described by its center (first the median, then the mean), its variability (also called  
805 spread—described both qualitatively and via a numerical measure, either inter-quartile  
806 range (IQR), range, or mean absolute deviation), and an overall shape (including  
807 descriptors such as symmetric, skewed left or right, peak, gap, and outlier) (6.SP.A.2,  
808 6.SP.A.3). As students explore datasets, they can produce visual representations of the  
809 distribution of their data; they can look at the shape of distributions that have different  
810 measures of center and spread, and develop visual understandings of the shape of  
811 distributions.

812 Students should have experiences, beginning in sixth grade, deciding which measure of  
813 center is a more useful descriptor of a typical value for data sets with different shapes.  
814 Because the mean is sensitive to extreme values, the median is often a more useful  
815 measure for skewed distributions; in this case, the inter quartile range is a useful  
816 measure of variability. For some distributions—with multiple clusters, for example—  
817 students may decide that neither median nor mean is a useful measure, and might  
818 decide that a single number cannot reasonably represent a typical value (6.SP.B.5).

819 Two tasks that reinforce the notion of these standard measures and replace rote  
820 disconnected calculation with conceptual thought are the following:

821 A. Students form a “name count line” creating a human graph to depict how many  
822 letters are in their first name. (All the students with five-letter names stand in a  
823 line, all those with four letters form a similar line to one side, and those with six  
824 letters form a line to the other, etc.) Then the teacher instructs one student from  
825 each end of the human graph to sit down. After repeating this multiple times, only  
826 one or two student(s) are still standing. If one, that student represents the median  
827 name length. If two, the median name length is halfway between the name  
828 lengths of the standing students.

829 B. Students are invited to explore the CODAP dataset of 4 elephant seals:  
830 [https://codap.concord.org/releases/latest/static/dg/en/cert/index.html?url=https://c](https://codap.concord.org/releases/latest/static/dg/en/cert/index.html?url=https://concord-consortium.github.io/codap-data/SampleDocs/Science/Biology/four-seals/Four_Seals.codap)  
831 [oncord-consortium.github.io/codap-data/SampleDocs/Science/Biology/four-](https://codap.concord.org/releases/latest/static/dg/en/cert/index.html?url=https://concord-consortium.github.io/codap-data/SampleDocs/Science/Biology/four-seals/Four_Seals.codap)  
832 [seals/Four\\_Seals.codap](https://codap.concord.org/releases/latest/static/dg/en/cert/index.html?url=https://concord-consortium.github.io/codap-data/SampleDocs/Science/Biology/four-seals/Four_Seals.codap)

833 The dataset includes data on the paths taken by the seals – visible on a mapping  
834 tool, the distance they swim, their latitude and longitude, the depth and  
835 temperature of the water and more.

836 In groups students are invited to explore the data and form investigative  
837 questions. Students start by plotting different variables with the graph tool, to  
838 consider the shapes of distributions. They choose to display the mean and  
839 median and consider how the measures of center relate to the visual distribution  
840 of the data. They form questions they are curious about: Do certain seals prefer  
841 deeper water? Does the distance seals swim relate to the temperature of the  
842 water? As students explore these questions they plot two variables on a graph  
843 and consider the slope of the relationship, they even add a third variable which is  
844 shown through color coding. Students learn to be comfortable investigating data,  
845 making use of measures to learn about their data.

846 Visual representations of distributions include box plots and histograms in sixth grade,  
847 adding to the line plots (called dot plots from grade six onward) from earlier grades

848 (6.SP.B.4). In addition, students learn to report and interpret measures of center and  
849 variability, and descriptions of distributions, in the context in which the data arose  
850 (6.SP.B.5). Seventh- and eighth-grade standards do not include additional  
851 representations of single-variable data sets, but these students should continue to  
852 create visual displays of such distributions.

853 In seventh grade, comparisons between two populations with similar variables is a  
854 context in which students describe and create visual displays of data. They can plot  
855 data and draw from different statistical methods such as creating box plots and dot plots  
856 to informally assess the degree of overlap of two populations, and students should be  
857 able to describe the difference between the two centers in terms of the measure of  
858 variability they use for the distributions.

## 859 Vignette

860 Alex did not enjoy learning about mean, median, and mode. He often confused the  
861 different measures and felt they had little meaning. His parent contacted Maria, his  
862 teacher, to let her know that he was expressing frustration about the meaning of the  
863 terms since his last assessment. Maria realized Alex was not alone since many of the  
864 students were still struggling with the meanings of these measures of average. Using  
865 this feedback as formative assessment, Maria approached the students with the idea to  
866 build physical models so they could experience the averages in visual and physical  
867 ways, encouraging important brain connections. Maria gave her students cubes and  
868 asked them to make 6 different towers of cubes that represented the numbers 1, 6, 3, 2,  
869 4 and 2. She asked them how they might construct a physical proof to show the mean  
870 of the numbers. Some of the students were able to calculate the answer; however, she  
871 kept pushing them to build a visual proof but remained open to multiple means of  
872 representation. This strategy, based on the UDL guidelines, allowed Maria to provide  
873 students scaffolds and supports to help highlight the patterns of language, and draw on  
874 background knowledge to express what they know in ways that are authentic and  
875 meaningful. Alex and his group members came up with the idea of moving the cubes  
876 form tower to tower to show that they could make six towers that were all the same

877 height. They just needed to average out all of the blocks. Alex and his group excitedly  
878 explained to the class how they had made a physical proof of finding the mean of the  
879 blocks. They shared the calculation with the class and compared it to the method they  
880 used of moving the blocks. After her students had discussed finding mean, Maria asked  
881 them to make a visual proof for the median and the mode.

882 Sampling to understand a population: randomness, bias, how many?  
883 (grades 7–8)

884 Prior to seventh grade, students' work with data has focused exclusively on using data  
885 to understand, describe, and compare the particular collection of objects or situations  
886 that were observed or measured. For example, to calculate the median highest  
887 temperature on school days in September as above, students would record the highest  
888 temperature on *each* school day.

889 Seventh grade includes the first introduction to *sampling*, the process of collecting data  
890 from a subset of a population in an attempt to understand or describe the whole  
891 population. This represents a big jump in sophistication from earlier work. Early  
892 experiences with sampling should first describe the measured variables for the sample  
893 (favorite lunch, number of minutes looking at screens, recorded for all students in the  
894 sample for one week), followed by team and class discussions about whether the  
895 description extends. For instance, if all students who come in to play basketball before  
896 school are asked to track their screen usage for the week, the class should discuss  
897 whether they expect the average of 862 minutes to be close to the average for everyone  
898 at school—and if not all students, then perhaps close to the average for some smaller  
899 definable group of students. Many similar discussions, with some obviously non-  
900 representative samples, help students understand the idea of a *random sample*.

901 If researchers decide to gather data from 40 members of the population, then their  
902 collection of 40 members is *random* if it is chosen in such a way that every possible  
903 subset of size 40 has an equal chance of being selected. It is important for students to  
904 have multiple experiences selecting samples from known populations in ways that are  
905 random (for instance, drawing numbered ping-pong balls from an opaque bag or

906 drawing student names on identical slips of paper from a hat) *and* in ways that are not  
907 random (for instance, asking survey questions only of the students who sit near you in  
908 class). The goal is an understanding that random sampling tends to produce samples  
909 that are *representative* of the population—that is, their distribution of the quantities  
910 under consideration are close to the distribution for the population as a whole  
911 (7.SP.A.1)—and a sense for the variability when using samples to make inferences and  
912 estimates for a population (7.SP.A.2).

913 Non-random sampling (such as attempting to understand the school as a whole by  
914 collecting data only from one’s friends, or by asking about eating habits at the gym after  
915 school, produce *biased* conclusions, even when the bias in the sample selection might  
916 not be obviously linked to the quantity being measured in the measurement or  
917 observation. *Bias* does not here refer to temperament or outlook (prejudice), which is  
918 one meaning of the word; instead; it means a *systematic error*.

919 Once random sampling becomes an available tool, the pool of questions open to  
920 students’ inquiry expands greatly: “I wonder how long on average it takes students from  
921 different grades to get from home to school?” “How do students who live in different  
922 areas spend time at the weekends?” “How much food is wasted in the lunchroom every  
923 month?” are all questions that could form a data exploration, as students consider their  
924 sample, which variables can be defined and collected, and engage in the four part  
925 exploration process.

926 Sampling is introduced in the seventh-grade standards and does not appear again until  
927 high school, but much of the eighth grade work with *bivariate* (two variable) data will  
928 make use of sampling, so it is important to continue activities that help understand  
929 *random sampling* through eighth grade as well. Students often believe that arbitrary  
930 sampling schemes (first 10 students I meet or every tenth student alphabetically) are  
931 random; they need to understand the difference between these schemes and choosing  
932 *by chance* so that every possible sample has an equal likelihood of being selected.



933 Vignette

934 Rosa has reflected on her seventh-grade students and how they have responded to the  
935 probability activities offered in her curriculum material. Overall, Rosa has not been  
936 satisfied with student understanding of random sampling. She decides to give students  
937 another, more visual and physical experience of the concept. Her plan calls for six  
938 paper bags filled with differently colored cubes. The sum of cubes and the color  
939 distribution of the cubes in the bags reflect the following:

940 Bag 1, 15 total: 15 blue

941 Bag 2, 12 total: 11 blue and 1 red

942 Bag 3, 20 total: 15 blue, 4 yellow, 1 red

943 Bag 4, 10 total: 5 red and 5 yellow

944 Bag 5, 12 total: 5 blue, 4 red, 3 yellow

945 Bag 6, 20 total: 8 blue, 8 red, 4 yellow.

946 Rosa explained to her students that their task was to determine the contents of each  
947 bag through sampling. She chose not to tell them how many times to sample but she  
948 did tell them to sample from the bags by selecting one cube at a time and then putting it  
949 back into the bag. Rosa also asked students to determine the chance of drawing a blue  
950 cube from each bag.

951 Students engaged in the activity, organizing how they would collect and record their  
952 information. When each group of students felt they had determined the number of cubes  
953 and color distributions of the contents of each bag, she asked them to choose which  
954 bag belonged to which card showing the contents of each bag. Rosa had filled the bags  
955 differently and made sure to have two different bags where the probability of drawing a  
956 blue cube would be one and another would be zero. After the activity and class  
957 discussion, Rosa was happy to hear her students later, talking about situations where  
958 the probability was one or zero as well as everything in between. Her students  
959 recognized the number of times they sampled usually led to better predictions about the

960 contents of the bags. They also realized sampling without replacement would have  
961 shown them the exact contents of the bag. The class engaged in a vibrant conversation  
962 about sampling with and without replacement, recognizing that it would be unproductive  
963 to draw all the cubes if there were a million.

#### 964 Are they related? Two changing quantities (grade 8)

965 Prior to grade seven, students work with a single collection of data measuring a single  
966 variable. In grade seven, they compare the same variable measured across two  
967 populations, either by actually measuring the whole populations or obtaining estimates  
968 for the distributions via sampling.

969 In eighth grade, the focus is *bivariate data*: Two quantities or categorical variables  
970 measured or observed across a population, or across a sample drawn from a  
971 population (8.SP.A.1). This work has important connections with linear equations and  
972 modeling.

973 The *scatter plot* as a visual representation of *quantitative* bivariate data is one of the  
974 most important ideas introduced here. A survey of students collecting both time and  
975 distance for traveling from home to school might reveal *clusters*, *outliers*, and any of  
976 various types of *association* (positive, negative, linear, non-linear). Students should  
977 describe such patterns in a scatter plot and interpret them in the context of the data  
978 (8.SP.A.1).

979 Students can explore large datasets—such as earthquake data from California—and  
980 explore bivariate relationships like, how does the location of earthquakes in the  
981 database compare to the magnitude of the earthquakes? They can plot the data using  
982 graphing tools and consider associations, data distributions, and relationships.

983 If students vary the weight added to a simple cart and measure the distance it travels  
984 when released at the top of a ramp, then plot the results on a distance (vertical axis) vs  
985 added weight (horizontal axis), they will likely see a relationship. This association  
986 between the two variables can then be *modeled* by a line if the association appears  
987 roughly linear (line-shaped). In eighth grade, students choose a line to fit the data by

988 visual approximation on the scatter plot, and compare and argue for whose line fits  
989 “best” (8.SP.A.2). They then interpret the meaning of the slope and intercept of their  
990 chosen model line, and use the line to make predictions for one variable when the other  
991 variable is specified (8.SP.A.3)

992 Finally, eighth grade students use two-way frequency tables as a tool to see  
993 associations in bivariate *categorical* data (8.SP.A.4). For instance, they might survey  
994 their class, including questions about favorite color and favorite genre of books, then  
995 input the data into a spreadsheet, organize the data and calculate relative frequencies  
996 in rows to explore possible relationships between the two variables.

### 997 What are the chances? Probability as the basis for data-based claims

998 Randomly selecting from a population and measuring a characteristic (in which variation  
999 is expected across the population) is a *chance process*: It may result in different results  
1000 and its outcomes follow some *distribution*.

1001 Probability is the way we express the chance of an outcome as a number between 0  
1002 and 1 (7.SP.C.5). Probability is combined with statistics in the grade seven standards;  
1003 statistics and probability are historically linked because statistical claims and estimates  
1004 are based on the mathematical field of probability. Models that draw from data science  
1005 and offer predictions of events, such as voting in elections, draw from probabilistic  
1006 reasoning.

1007 The connections between probability and statistics is often not clear to students,  
1008 especially when their experiences focus on procedures and calculation rather than  
1009 exploration, context and interpretation. There is much work with probability that does not  
1010 support statistical reasoning (for example, calculating theoretical probabilities for the  
1011 sum of two dice without using those theoretical probabilities to decide whether a given  
1012 pair of dice are likely fair), and middle school probability experiences should be carefully  
1013 designed to support reasoning with interesting and meaningful data.

1014 In seventh grade, students gather data to estimate the probability of outcomes by  
1015 observing their long-run relative frequency; that is, they compute *experimental*

1016 *probability*. Consider repeating this experiment 150 times: draw a marble from a bag  
1017 with marbles in it, record its color, then put the marble back in the bag. If we get a blue  
1018 marble 32 times, our estimate for the probability of getting blue on any particular draw is  
1019  $32/150$  (7.SP.C.6, 7.SP.C.7.B).

1020 Compare the marble experiment just described to this one: Put the following marbles in  
1021 a bag (all identical except for color): 16 blue marbles, 31 red marbles, 16 green  
1022 marbles, and 12 white marbles (75 total marbles). If you blindly pull a marble from the  
1023 bag, what is the probability that you will get a blue marble? If you do this 150 times  
1024 (putting the marble back each time), about how many times do you expect to get a blue  
1025 marble? After calculating this expectation, students should construct an algorithm or  
1026 pseudo-code to run the simulation 150 or 1500 or 15,000 times to compare with their  
1027 theoretical expectations (CS Standards 6-8.AP.10).

1028 Note the difference between the questions in the previous two paragraphs: In the first,  
1029 students use long-run relative frequency to estimate probability; in the second, students  
1030 build a (*theoretical*) probability model and use it to estimate long-run frequency  
1031 (7.SP.C.7). If a marble experiment is then performed and relative frequencies of  
1032 outcomes do not seem close to predictions from the probability model, then students  
1033 need to be able to discuss possible sources of discrepancy (7.SP.C.7): Perhaps the  
1034 green marbles feel different and tend to be drawn more frequently than predicted.  
1035 Maybe somebody changed the mix of marbles in the bag. Or perhaps not enough draws  
1036 were performed to see the relative frequencies approach the probability model.

1037 Finally, seventh grade students find probabilities of compound events (events which are  
1038 made up of several simple events; for example, drawing two marbles from the bag of 75  
1039 described above and getting one white and one blue marble) (7.SP.C.8).

1040 The specific calculations above are not central to the data science progression, but  
1041 recognition that some events (repeat the draw 5 times, get all blue; or repeat the draw 5  
1042 times, obtain WBWWB in that order) are *much* less likely than others (repeat the draw 5  
1043 times, get 3 white and 2 blue) is key to understanding claims made from data.

1044 In fact, most statistical claims depend on a comparison of a (theoretical and  
1045 hypothetical) probability model with observed data, as in 7.SP.C.7. For middle school  
1046 students to be prepared for future data science work, they need experiences to build a  
1047 sense that more data tends to produce relative frequencies closer to actual probabilities.

1048 Invite students to explore rich datasets, such as the distribution of births in the US – and  
1049 consider questions of probability, that they can explore, such as: what is the chance  
1050 that two people share the same birthday? This is a question that could be explored  
1051 theoretically or experimentally. (More at  
1052 [https://codap.concord.org/releases/latest/static/dg/en/cert/index.html?url=https://concord-  
-consortium.github.io/codap-  
data/SampleDocs/Mathematics/Probability/Birthdays/Birthdays.codap](https://codap.concord.org/releases/latest/static/dg/en/cert/index.html?url=https://concord-<br/>1053 -consortium.github.io/codap-<br/>1054 data/SampleDocs/Mathematics/Probability/Birthdays/Birthdays.codap))

## 1055 Vignette

1056 Quincey started middle school without a lot of interest in math class. Quincey had  
1057 always been interested in how the world works, and science and social studies were  
1058 their favorite classes. Quincey had not had much experience of math class connecting  
1059 to areas of interest.

1060 Quincey’s sixth-grade math teacher, Lori, was determined to change this for all  
1061 students. Lori knew that the data science standards in sixth grade would give Quincey  
1062 an opportunity to use real data from the world and to understand that they could ask  
1063 questions of data and see the connections between mathematics and life. Lori decided  
1064 to use an activity to explore the “shape” of data: The context is hurricanes in the Atlantic  
1065 Ocean using real data collected from five years of hurricanes spread over four decades.  
1066 Quincey loved the opening discussion in class when the students all discussed the 2017  
1067 hurricane data displayed on a line plot. Quincey and the class were really interested in  
1068 the number of hurricanes that were in category 0—tropical storms. Next, students  
1069 worked in groups where they studied hurricane category data for the years 1977, 1987,  
1070 1997 and 2007. Each decade’s data was presented in different ways: bar graph, line  
1071 plot, tables and sentences. Quincey enjoyed the analysis and was taken with the  
1072 different ways of displaying data as well as the changes in the spread of data.

1073 Throughout the discussions Quincey asked questions about the science of hurricanes.  
1074 *How do they develop?* he wondered. *What makes them get larger? What is the*  
1075 *difference between a category 3 storm and a category 5 storm?* At the close of the  
1076 lesson Lori was convinced that students understood that different visual displays of data  
1077 can make it easier to see the shape of data. The shape of the data on the displays  
1078 helped students see how a situation might be changing over time. The class reflected  
1079 that the changes were easier to see in line plots and histograms versus the data being  
1080 shared in writing or in a table of values. Quincey decided to further investigate the  
1081 number of category 4 and 5 hurricanes over the past 100 years and how these storms  
1082 become stronger, and they set out to gather more data and ask questions of the data.  
1083 Others in the class decided to investigate why the number of category 4 and 5 storms  
1084 are increasing.

## 1085 High School

1086 In this outline for data science in high school, two sections of guidance are provided: (1)  
1087 experiences and expertise in data science common for all high school students, and (2)  
1088 experiences and expertise for a high school pathway with a data science focus  
1089 (expanding on the pathway outline in Chapter 8, Grades 9–12).

1090 Computers have become central to modern life. Whether laptops, phones, so-called  
1091 “smart” appliances, medical records systems, exercise trackers, GPS location  
1092 recording, payment methods, or other forms, computers are involved in most of our  
1093 transactions. Every interaction with a computer generates data about that interaction  
1094 (which is collected and saved)—but very little of this data is analyzed and interpreted.

1095 Even as computers have led to the collection of vast amounts of data, computational  
1096 tools (including both computer hardware capabilities and advances in algorithms) have  
1097 dramatically altered the available methods for making use of and communicating  
1098 interpretations of data. In fact, meaningful analysis of large or multivariate data sets is  
1099 impossible without computer tools.

1100 For many questions about which students might wonder, existing data sources might  
1101 provide the necessary information. Designing data collection to obtain exactly the

1102 desired data for answering a specific question (the classical statistical experiment  
1103 approach, still the main approach in grades K–8 above) is expanded to include  
1104 descriptive techniques for analyzing multivariate data, critical questioning skills to  
1105 interrogate pre-existing data's suitability for the investigation, and an understanding of  
1106 ways to access and acquire data through the internet. This understanding-extraction  
1107 uses two processes: data description methods, both visual and numerical, to investigate  
1108 conjectures and discover patterns; and model-building to test conjectures, make  
1109 predictions of future observations, and evaluate the predictive success. These huge  
1110 existing data sets are not collected in order to answer a particular question, do not  
1111 typically represent random samples, and are often missing data or are otherwise  
1112 “messy.”

1113 Data science should be understood as a broad term encompassing many tools relevant  
1114 to learning from data. These include tools of traditional statistics classes, but also  
1115 include computational tools to address the massive size and complexity of many of  
1116 today’s data sets, and disciplinary knowledge of the field generating the data. Thus,  
1117 data science is an inherently interdisciplinary field that uses scientific and statistical  
1118 methods and processes to derive understanding, insight, and predictive ability from  
1119 (often unstructured) data (Dhar, 2013).

## 1120 Data science for equity and inclusion

1121 An important way in which educators can offer social and emotional support to students  
1122 is by designing engaging lessons that allow students to connect with the ideas being  
1123 taught. Traditional mathematics lessons that have taught mathematics as a set of  
1124 procedures to follow have resulted in widespread disengagement as students see no  
1125 relevance for their lives. This is particularly harmful for students of color and for girls—  
1126 who receive additional harmful messages that mathematics is not for them. Data  
1127 science is a field that provides opportunity for equitable practice, with multiple  
1128 opportunities for students to connect with ideas and for them to receive messages that  
1129 data science is a field in which any students can excel.

1130 Walton and colleagues (2015) have shown through numerous studies that many  
1131 students, particularly girls and students of color, do not feel that they belong in certain  
1132 disciplines. This is often due to a history of negative and off-putting messages (Chestnut  
1133 et al, 2018). Other studies have shown that different topics and teaching approaches  
1134 can lead to feelings of belonging or not belonging (Boaler, 2019; Boaler, Cordero &  
1135 Dieckmann, 2019). Data science is well placed for teachers to offer students a sense of  
1136 belonging as students are invited to investigate real data that will often be relevant to  
1137 their lives. This meaningful engagement will offer students opportunities to develop self-  
1138 confidence and self-efficacy. When teachers offer students the role of being data  
1139 investigators, asking and answering questions with data, they can take an active role in  
1140 their learning, encouraging opportunities for self-motivation and goal setting. Important  
1141 principles in the teaching of data science, that will offer the greatest chance for social,  
1142 emotional, and academic development, include the following:

1143 ● *Mindset and Belonging Messages*

1144 At frequent times students should be reminded that data science is a field in  
1145 which all people are welcome and can succeed. In line with successful  
1146 interventions in mindset and belonging students should also be reminded that  
1147 struggle is an important part of learning, that all students have times of struggle,  
1148 and that the difference between successful and unsuccessful students is the way  
1149 they respond to times of difficulty. Share with students examples of successful  
1150 people inside data science, that highlight gender and racial diversity (see for  
1151 example:

1152 [https://www.youtube.com/watch?v=KYvhoH5AzHA&feature=emb\\_logo](https://www.youtube.com/watch?v=KYvhoH5AzHA&feature=emb_logo)).

1153 ● *Use Real Data*

1154 Data science gives an opportunity for students to be asking questions of real  
1155 data sets, developing social awareness, and investment in the solutions they  
1156 discover. When working with secondary data sets (data obtained from others,  
1157 rather than collected by students), teachers should choose meaningful ones in  
1158 order to give students an important connection to the content they are learning  
1159 and opportunities to take the perspective of others, which will help them develop  
1160 empathy. When teachers use local data sets they can also help students feel like



1161 they are important members of their community – as they explore questions and  
1162 find answers to local problems that they can help with real data. Identifying  
1163 problems and finding solutions will help students develop responsible decision-  
1164 making.

1165 Some teachers worry that they cannot provide culturally sustaining connections  
1166 for their classes as they are not experts in the cultures of all their students, but  
1167 real data sets from different communities provide opportunities for students to  
1168 bring their own knowledge and expertise to data rich problems. There should  
1169 also be times when students are invited to collect data from their own community  
1170 and build their own data sets. Students can pose questions that are important to  
1171 them, including those with cultural meaning, collecting data from their own lives  
1172 and communities. As Paris (2012) describes students will be fostering and  
1173 sustaining “linguistic, literate, and cultural pluralism.” The act of collecting data  
1174 provides an important learning opportunity for students to understand decisions  
1175 that need to be made around the collection and organization of data as well as  
1176 how to deal with uncertainty in their data. Students will be the ones with  
1177 important expertise in these investigations.

1178 ● *Focus on Collaboration and Communication*

1179 Data science is a field in which people collaborate, connecting ideas to solve  
1180 difficult problems with data. Ideal collaborations are those in which diverse  
1181 groups of students come together and work effectively with different ideas being  
1182 valued and developed. Groupwork is enhanced when students are given open  
1183 problems in an environment where different students feel safe to share their  
1184 ideas, respectfully, before working to solutions. Group work is usually much more  
1185 effective when students start their work in structured and unstructured  
1186 conversations where each group member shares their thoughts. Collaborative  
1187 classrooms founded in engaged listening and the capacity to articulate verbally  
1188 as they build on each other’s ideas, are places where students feel valued and  
1189 where they develop Important relationship skills of communication, social  
1190 engagement and teamwork.

1191 Data for all: living in an information-overloaded world

1192 Because decisions and predictions are often based on data, all California high school  
1193 graduates need data acumen, as described in the chapter introduction: skills in  
1194 interpreting and visualizing data, making and critiquing data-based arguments, and  
1195 some facility with data software. The ability to identify types of questions that are subject  
1196 to exploration through data is crucial, as is an understanding of some misuses of data  
1197 and of one’s own online data footprint. As in earlier grades, it is crucial that students  
1198 have opportunities to generate and investigate their own “I wonder” questions in given  
1199 contexts. All statistics standards are identified as *modeling* standards, reflecting the  
1200 origin of all work with data in authentic questions about the world.

1201 “I wonder” is just the beginning, of course. Students must learn to formulate statistical  
1202 investigative questions, pose data collection questions, interrogate existing data,  
1203 analyze data, and formulate, interpret, and communicate findings. Note that questioning  
1204 is a central practice throughout the statistical problem-solving process. As the GAISE II  
1205 report summarizes:

1206       The statistical problem-solving process typically starts with a statistical  
1207       investigative question, followed by a study designed to collect data that aligns  
1208       with answering the question. Analysis of the data is also guided by questioning.  
1209       Constant questioning and interrogation of the data throughout the statistical  
1210       problem-solving process can lead to the posing of new statistical investigative  
1211       questions.

1212       Often when considering secondary data, the data need to first be interrogated—  
1213       how were measurements made, what type of data were selected, what is the  
1214       meaning of the data, and what was the study design to collect the data. Once a  
1215       better understanding of the data has been gained, then one can judge whether  
1216       the data set is appropriate for exploring the original statistical investigative  
1217       question or one can pose statistical investigative questions that can be explored  
1218       with the secondary data set. (Bargagliotti et al., 2020)

1219 The process of *interrogating* secondary data (data collected by anyone other than the  
1220 person doing the analyst) is crucial. Public datasets are not collected specifically to  
1221 answer students' questions. So before using such datasets, students will need to  
1222 evaluate the appropriateness for their purposes: How do the measures, methods, and  
1223 scope of the dataset match the statistical investigative question(s) that interest us? For  
1224 what purpose(s) were the data collected?

1225 Using data to answer authentic questions about the world is a powerful antidote to the  
1226 famous student retort "when will I ever need to know this?" The Mathematics:  
1227 Investigating and Connecting high school pathway in Chapter 8: Grades 9–12,  
1228 encourages the use of data science contexts to frame many of students' explorations to  
1229 develop content and practice standards in all domains. In this section, we propose an  
1230 outline for data science understanding, considering the understandings high school  
1231 students should develop.

1232 Students enter high school with significant relevant experiences which should be drawn  
1233 upon in high school work. As the CA CCSSM introduction to High School Probability  
1234 and Statistics summarizes work in prior grades:

1235 Data are gathered, displayed, summarized, examined, and interpreted to  
1236 discover patterns and deviations from patterns. Quantitative data can be  
1237 described in terms of key characteristics: measures of shape, center, and spread  
1238 [variability]. The shape of a data distribution might be described as symmetric,  
1239 skewed, flat, or bell shaped, and it might be summarized by a statistic measuring  
1240 center (such as mean or median) and a statistic measuring spread (such as  
1241 standard deviation or interquartile range). Different distributions can be compared  
1242 numerically using these statistics or compared visually using plots. Knowledge of  
1243 center and spread are not enough to describe a distribution. Which statistics to  
1244 compare, which plots to use, and what the results of a comparison might mean,  
1245 depend on the question to be investigated and the real-life actions to be taken.

1246 The big ideas of data science for all students in high school are identified in the  
1247 statistics cluster headings in the standards, with an additional big idea discussed here,

1248 in response to the changing approaches to data described above. The first two are  
1249 described in more detail, with additional examples, in the Draft High School Progression  
1250 on Statistics and Probability (<https://www.math.arizona.edu/~ime/progressions/>).

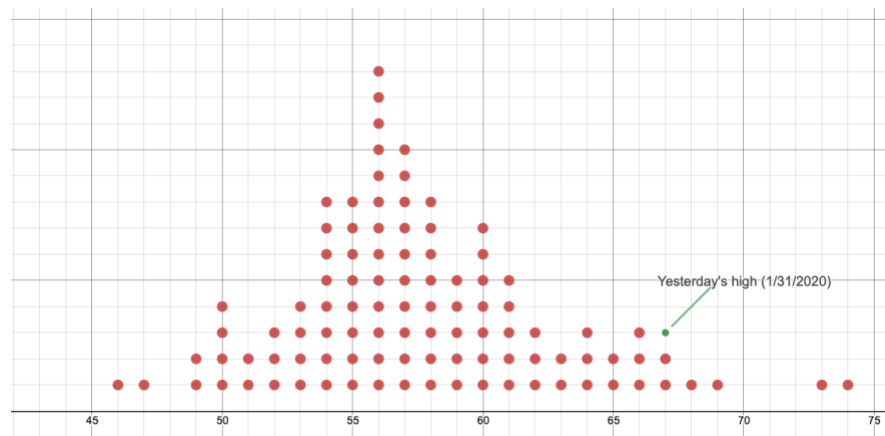
- 1251 • Interpreting Categorical and Quantitative Data
- 1252 • Making Inferences and Justifying Conclusions
- 1253 • From statistics to data science

#### 1254 Interpreting Categorical and Quantitative Data

1255 High school students continue their work with the representations of data introduced in  
1256 K–8, and they are introduced to the normal curve as an approximating model for some  
1257 data sets. However, the major work in high school in interpreting data is using functions  
1258 as models of associations in two-variable quantitative data.

1259 Building on K–8 experiences, high school students continue to visualize and represent  
1260 *single-variable* data with dot plots, histograms, and box plots; use measures of center  
1261 and spread to describe such distributions; and compare distributions from different  
1262 populations or samples using these representations and statistics (S-ID.1–3). For data  
1263 sets that appear to be bell-shaped, they use the mean and standard deviation to specify  
1264 an approximating normal distribution and to approximate population percentages in

1265 specified ranges (S-  
1266 ID.4). Students might,  
1267 for example, obtain  
1268 high temperatures on  
1269 a specific date over  
1270 the past 100 years  
1271 from a nearby  
1272 weather station



1273 (<https://calclim.dri.edu/pages/stationmap.html>), create a dot plot, visually check for a  
1274 bell-shaped distribution, and use an approximating normal distribution to make a case  
1275 for whether or not the temperature was consistent with historical trends (CA  
1276 Environmental Principles & Concepts 1.C).

1277 When available data includes two (or more) measurements or characteristics for each  
1278 observation, students' tools for representing and interpreting relationships depend on  
1279 the nature of each variable.

- 1280 ● If both are categorical, two-way frequency tables give an important summary that  
1281 reveals relationships when interpreted in the context of the data
- 1282 ● If one is categorical and the other quantitative, students can treat each category  
1283 as a separate population and compare the quantitative data for the different  
1284 categories as in the single-variable paragraph above
- 1285 ● If both variables are quantitative, the scatter plot is the standard visual  
1286 representation

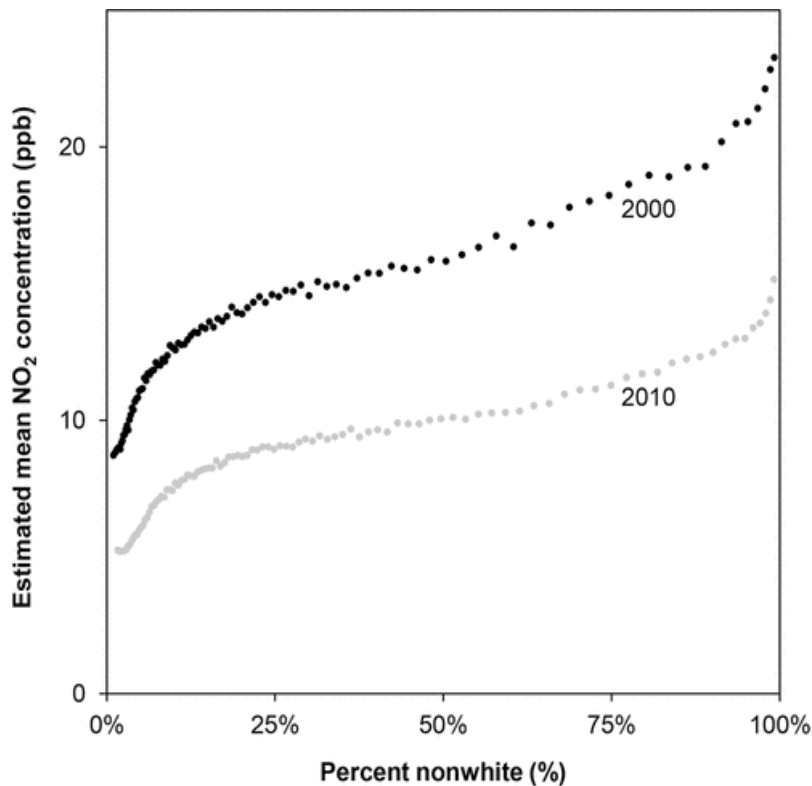
1287 **Modeling association in numerical data using functions:** Once a scatter plot is  
1288 created, an association between the two variables may become visually identifiable.  
1289 Fitting a function to the data is the creation of a mathematical model for the association.  
1290 This begins in eighth grade with visual fitting of a linear model. While the type of  
1291 function that is used most frequently is a line (a linear function), students also need  
1292 experiences with plotting associations that are clearly non-linear, as well as experiment  
1293 with fitting other types of functions (quadratic, exponential).

1294 Any standard data software (including spreadsheets, Desmos, Geogebra, CODAP) will  
1295 fit lines, quadratic functions, and exponential functions to given data. The specific  
1296 standard technique for identifying a line (or quadratic or exponential function) of best fit  
1297 (least-squares regression) is *not* an expectation; but students should have experiences  
1298 fitting lines and some other functions visually (by adjusting parameters on appropriate

1299 function types in graphing software) and using appropriate software tools which perform  
1300 the regression behind the scenes.

1301 Most importantly, functions which model associations must be used to solve problems  
1302 (e.g. prediction of one variable given another) (S-ID.6.a) and must be interpreted in the  
1303 context of the data (S-ID.7).

1304 Important examples of modeling association in numerical data arise in many contexts in



1305 science, history, physical education, and social studies. The California History–Social  
1306 Science Framework (in its Appendix C) expects students to develop *Chronological and*  
1307 *Spatial Thinking*, including analyzing change over time; both *time* and *space* provide  
1308 opportunities for finding meaningful quantities that vary together. For example, students  
1309 might wonder whether pollution exposure is related to wealth, and either find zip-code  
1310 level data on both air pollution and income, or find existing research like the graph here,  
1311 and work to understand and explain it. (This graph has the added benefit of  
1312 representing both change over time and change in space.)

1313 Source: <https://ehp.niehs.nih.gov/doi/10.1289/EHP959>.

1314 Making Inferences and Justifying Conclusions

1315 Making conclusions and generalizations about a population from a sample (S-IC.1) is  
1316 the goal of *inferential* statistics, as opposed to *descriptive* statistics. Students work with  
1317 random samples beginning in seventh grade, their first experience trying to understand  
1318 a population without gathering data about all of its members. This strand of high school  
1319 data work is the foundation for most meaningful use of statistics for making decisions.

1320 Students must decide whether a result observed through data is consistent with a  
1321 mathematical model of the process that generates the data (S-IC.2). For instance, if a  
1322 student estimates that 30 percent of the students at the school grow food at home, then  
1323 that is a mathematical model that gives them an idea of what proportion to expect in a  
1324 sample. If they then survey 5 randomly-chosen students, and all say they grow food at  
1325 home, then the student should be able to reason as follows: *If* 30 percent of students  
1326 grow food at home, then the chances of five randomly-chosen students all being among  
1327 those 30 percent of students is  $(.3)^5 = .00243 = .342$  percent, or less than half of one  
1328 percent. Thus, the student might doubt—that is, they might *reject*—the 30-percent  
1329 hypothesis. Students should have many experiences of simple situations like this to  
1330 understand how decisions based on data rely on probability, and are not *guaranteed* to  
1331 produce correct answers to the original question.

1332 Students should work with data originating in four different methods of data production,  
1333 including at least some student-generated questions and student-gathered data. These  
1334 methods are (1) *census* data; that is, data that contain measurements on every member  
1335 of the target population (such as the database of crimes occurring in a given city in a  
1336 given time frame, or rain gauge data for a given location, which captures all precipitation  
1337 at that location); (2) surveys administered to random samples (to estimate population  
1338 values, or *parameters*, for the surveyed quantities), (3) randomized experiments (to  
1339 compare treatments and demonstrate cause), and (4) observational studies (to study  
1340 characteristics or quantities when random selection or assignment is not possible) (S-  
1341 IC.3). The Draft High School Progression on Statistics and Probability  
1342 (<https://www.math.arizona.edu/~ime/progressions/>) contains detailed examples  
1343 describing the treatment of each that is expected in the standards.

1344 For surveys and experiments, the key understanding is the link between the random  
1345 selection or assignment and the ability to reason probabilistically to make claims. With a  
1346 survey, the random sampling allows generalizing to a population. With an experiment,  
1347 the random assignment allows causal conclusions but not generalization to a broader  
1348 population—unless the sample in the experiment was randomly selected from some  
1349 larger population. For example, medical studies (experiments) must use willing  
1350 volunteers and thus are not random samples of the overall population; this makes it  
1351 much harder to draw broadly-applicable conclusions.

1352 In a college statistics course or the data science course outlined below, students will  
1353 learn ways to quantify the comparisons between gathered data and hypothesized  
1354 population parameters (margins of error and  $p$ -values). Making sense of these,  
1355 however, requires an understanding of the role of randomization in the data gathering.

1356 When using a sample mean or proportion to estimate a population mean or proportion,  
1357 students use simulation models to estimate a margin of error, instead of formulaic  
1358 calculations. Briefly, the process is to use data simulation software to draw many  
1359 random samples from a hypothetical population, and to see how often a result is  
1360 obtained that is as extreme as the sample mean or proportion. Doing this process for  
1361 hypothetical populations with many different mean or proportion parameters helps  
1362 students see that there is a range of population parameters that often (more than 5% of  
1363 the time) produce simulated sample means or proportions that are as extreme as (or  
1364 more extreme than) the actual sample mean or proportion. This range of population  
1365 parameters is the (simulation-based) confidence interval, given as (sample mean or  
1366 proportion  $\pm$  margin of error). Note the probabilistic argument here: *If* the population  
1367 mean or proportion were outside of the confidence interval, *then* sample means or  
1368 proportions as extreme as we obtained in our random sample would be rare. So, we  
1369 expect that the true population mean or proportion is within the confidence interval. (*But*  
1370 *cannot be certain* that it is!)

1371 A similar process is used to evaluate confidence in a randomized experiment, in which  
1372 subjects are randomly assigned to two or more treatment groups. (Treatment could  
1373 mean medical treatment, or assignment of different tasks, or being shown different



1374 motivational videos, etc.) Some quantity is then measured for each subject, and the  
1375 investigator then has to decide from the results whether a treatment, say treatment A,  
1376 produced any effect on the measured quantity. Simply having a different mean for each  
1377 treatment group is not enough, as we expect variation in the measurement and thus  
1378 between groups. In this case, all of the treatment groups are pooled into a population  
1379 and then re-sampled (randomly) many times, to see how often the re-sampled mean or  
1380 proportion is at least as extreme as the actual treatment A group difference. If such  
1381 differences are rare, the experiment is taken as evidence that treatment A caused a  
1382 change in the measured quantity.

### 1383 From Statistics to Data Science

1384 For questions about community, society, or natural systems beyond students'  
1385 immediate experience—and thus beyond their ability to gather data directly—existing  
1386 data can often be identified from online sources. Since students frequently encounter  
1387 claims made from such large data sets, it is crucial that all students have experiences in  
1388 which they explore the ways in which such claims are made. A major difference  
1389 between the classical statistical approach begun in K–8 and the “big data” of the  
1390 growing field of data science is the richness and complexity of available data sets, even  
1391 more so than their sheer size.

1392 Many sophisticated approaches to working with rich, complex data sets are left to a data  
1393 science course in the data science pathway; but *all* high school students should  
1394 exercise and refine their understanding of data exploration, causal inference, and  
1395 statistical reasoning using large, real world data sets. As students work with these data  
1396 sets they can draw upon the data science understandings they have developed in their  
1397 K–8 mathematics lessons. Emphasis should be on questioning and interpreting, rather  
1398 than technical procedures.

1399 Data exploration begins with a search for available data about a context of interest. The  
1400 data set is then examined for hidden patterns and associations (usually via visual  
1401 representations). Any patterns or associations discovered can lead to new hypotheses  
1402 or questions to investigate further. Students began this process in eighth grade, and

1403 continue in high school with experiences in which they examine data sets with multiple  
1404 variables measured for each member of the sample. They plot pairs of variables to  
1405 decide which ones might show associations. Important discussions for students to  
1406 engage in when working with existing data sets include

- 1407 ● Prior to exploring: Do we *expect* any of these variables to be associated? Why?
- 1408 ● Might the association we see just be a result of the way in which the data was  
1409 collected, rather than truly reflective of the population? What features of the data  
1410 collection might make conclusions suspect, and what features might give  
1411 confidence? Note that a large sample size is not enough to have confidence in  
1412 conclusions.
- 1413 ● Can we think of possible explanations for the association(s) we see? Can we  
1414 think of ways we could decide which explanations might be accurate?

1415 After data exploration identifies some association(s) of interest, the stage of model  
1416 building follows. Technical methods are reserved for the specialized data science  
1417 course below, but all students need to explore questions such as:

- 1418 ● Could we use some variables to predict others? This is a hugely important use of  
1419 data, since some factors are easier to measure or observe than others. In  
1420 medicine and many other fields, this often takes the form of trying to predict  
1421 future outcomes using presently-available information.
- 1422 ● If we could only know measure one variable to try to predict a variable of interest,  
1423 which one would we pick? Why? What if we could measure two? Which second  
1424 variable gives us the most *new* information for prediction?

1425 Most importantly, high school students (like K–8 students) must experience data  
1426 science as a set of tools for making sense of their worlds in ways that matter to them.

### 1427 Advanced high school data science

1428 The traditional sequence of high school courses—algebra, geometry, algebra 2—was  
1429 standardized in the United States following the “Committee of Ten” reports in the 1890s.

1430 The course sequence—which was primarily designed to give students a foundation for

1431 calculus—has seen little change since the Space Race in the 1960s. With the rapid  
1432 expansion of information available to all in the form of data, far more students pursue  
1433 statistics classes than calculus, and may be better served by a data science course as a  
1434 culminating high school mathematical science experience. In addition to the importance  
1435 of the data science content—to 21st Century jobs and to a wide range of college  
1436 majors—many students are more engaged by open-ended explorations of important  
1437 data sets, drawing upon important mathematical principles and tools, than by many  
1438 traditional courses organized around mathematical techniques. This Framework  
1439 provides design principles and content outcomes for such a course.

1440 California high schools offer upper-level data science courses in two ways. In the first  
1441 model, students have a common experience in grades nine and ten, with pathways  
1442 branching at grade eleventh. Some districts have designed and are offering eleventh  
1443 grade data science courses as an option for this third year of high school mathematics;  
1444 in this case, the ninth and tenth grade courses need to be designed to include the  
1445 important high school geometry standards. The second model is a data science course  
1446 as a fourth-year course, following a coherent three-year pathway that builds the “for all  
1447 students” data science understanding outlined in the previous section. The design  
1448 principles and content outcomes below are flexible enough to be implemented in either  
1449 model, with appropriate adjustments for students’ prior experiences.

#### 1450 Design Principles

1451 These principles provide guidelines for design of curricular materials and classroom  
1452 instruction for a data science course, in order to support a coherent and engaging  
1453 experience for students. These principles should be used by developers to build  
1454 curricular materials that are true to the vision of the course, as well as by educators  
1455 reviewing materials and developing a repertoire of pedagogical strategies for use in  
1456 teaching the course. Many students and teachers already engage in these behaviors; in  
1457 these cases, these design principles will be seen as reinforcing and supportive. The  
1458 spirit of this framework recognizes that, at some levels, everyone is a learner, and  
1459 everyone is growing an understanding of mathematics, each other, and the world we  
1460 share.

Design Principle	Students will	Teachers will
<p><b>Active Learning.</b> The course provides regular opportunities for students to actively engage in data explorations using a variety of different instructional strategies (e.g., hands-on and technology-based activities, small group collaborative work, facilitated student discourse, interactive lectures).</p>	<ul style="list-style-type: none"> <li>● Be active and engaged participants in discussion, in working on data explorations with classmates, and in making decisions about the direction of instruction based on their work.</li> <li>● Actively support one another's learning.</li> <li>● Discuss results of their explorations with the instructor and/or classmates in class.</li> </ul>	<ul style="list-style-type: none"> <li>● Provide low-floor, high-ceiling activities and explorations that all students can access and that extend to high levels. Such activities should provide meaningful opportunities for exploration and co-creation of mathematical understanding and data literacy.</li> <li>● Provide interesting and sometimes local data sets and invite students to ask questions of the data. Encourage different students to pose and investigate different questions, and to come together to discuss findings.</li> <li>● Facilitate students' active learning of data science through a variety of instructional strategies, including inquiry, problem solving, critical thinking, and reflection.</li> <li>● Create a safe, student-driven classroom environment in which all students feel a sense of belonging to the class and the discipline, are encouraged to take risks and embrace mistakes, and are able to make decisions about the direction for</li> </ul>

		<p>instruction through the results of their exploration of data science. Students' ideas are at the center of the conversation.</p>
--	--	-------------------------------------------------------------------------------------------------------------------------------------

<p><b>Growth Mindset.</b> Courses support students in developing the tenacity, persistence, and perseverance necessary for learning data science, for using mathematics and statistics to tackle authentic problems, and for being successful in post-high school endeavors.</p>	<ul style="list-style-type: none"> <li>● Make sense of data explorations by drawing on and making connections with their prior understanding and ideas.</li> <li>● Persevere in solving problems and realize that it is acceptable to say, “I don’t know what to do next,” but that it is not acceptable to give up</li> <li>● Seek help from different sources to move forward in their investigations.</li> <li>● Compassionately help one another by sharing strategies and solution paths rather than simply giving answers.</li> <li>● Reflect on mistakes and misconceptions to improve their mathematical understanding and data literacy.</li> <li>● Understand that struggle is valuable for brain growth and times of struggle should be valued.</li> <li>● Develop/strengthen a growth mindset to continue to apply in mathematics, data science, and other areas of their post-high-school life.</li> </ul>	<ul style="list-style-type: none"> <li>● Provide information about and model the importance of having a growth mindset.</li> <li>● Value mistakes and times of struggle.</li> <li>● Facilitate discussions on the value of mistakes, misconceptions, and struggles.</li> <li>● Give students time to struggle with tasks and ask questions that scaffold students’ thinking without stepping in to do the work for them.</li> <li>● Praise students for their efforts in making sense of mathematical ideas and for their perseverance in reasoning through problems and in overcoming setbacks and challenges in the course.</li> <li>● Provide students with low-stakes opportunities to fail and learn from failure.</li> <li>● Provide regular opportunities for students to self-monitor, evaluate, and reflect on their learning, both individually and with their peers.</li> </ul>
----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

<p><b>Problem Solving.</b> Courses provide opportunities for students to make sense of problems and persist in solving them.</p>	<ul style="list-style-type: none"> <li>● Apply intuition, life experience, and previously learned strategies to solve unfamiliar problems.</li> <li>● Explore and use multiple solution methods.</li> <li>● Share and discuss different solution pathways and methods.</li> <li>● Be willing to make and learn from mistakes in the problem-solving process.</li> <li>● Use tools and representations, as needed, to support their thinking and problem solving.</li> <li>● Develop and justify their own strategies to approach new problems.</li> </ul>	<ul style="list-style-type: none"> <li>● Present tasks that require students to find or develop a solution method.</li> <li>● Provide data sets that allow for multiple strategies and solution methods, including transfer of previously developed skills and strategies to new contexts.</li> <li>● Provide opportunities to share and discuss different solution methods.</li> <li>● Model the problem-solving process using various strategies.</li> <li>● Encourage and support students to explore and use a variety of approaches and strategies to make sense of and solve problems.</li> </ul>
<p><b>Authenticity.</b> Courses present data science as a subject and learning that allows us to model and solve problems that arise in the community.</p>	<ul style="list-style-type: none"> <li>● Recognize specific ways in which mathematics and data are used in everyday decision making.</li> <li>● Recognize problems that arise in the real world that can be solved with data science</li> <li>● Contribute meaningful questions that can be answered using data science.</li> <li>● Experience the decision making involved in collecting, cleaning,</li> </ul>	<ul style="list-style-type: none"> <li>● Provide opportunities to ask questions of data sets that are relevant to students, both in class and on assessments</li> <li>● Provide opportunities for students to pose questions that can be answered using data science methods and tools, and answer them.</li> <li>● Provide students with real data to explore and work with, including doing some of the data cleaning that is often required.</li> </ul>

	<p>analyzing, and visualizing data.</p>	
<p><b>Context and Interdisciplinary Connections.</b> Courses present data science in context and connects data science to various disciplines and everyday experiences.</p>	<ul style="list-style-type: none"> <li>● Contribute personal experiences, where appropriate, that connect to classroom experiences.</li> <li>● Actively seek connections between classroom experiences and the world outside of class.</li> <li>● Examine the ways that data is collected in their day-to-day lives, and consider the ethics and consequences of collecting and using data to make decisions.</li> </ul>	<ul style="list-style-type: none"> <li>● Provide opportunities for students to share their personal backgrounds and interests, including cultural values, and help make the connection between what is important in students' lives and future aspirations, and what they are learning in data science.</li> <li>● Provide real and interesting data sets, including some that are local to students</li> <li>● Invite students into data explorations that illustrate authentic applications.</li> <li>● Provide data explorations that include applications from a variety of academic disciplines, programs of study, and</li> </ul>



		careers, and which are culturally sustaining.
--	--	-----------------------------------------------

<p><b>Communication.</b> The course develops students' ability to communicate their data explorations and findings in varied ways including with words, data visualizations and numbers.</p>	<ul style="list-style-type: none"> <li>● Present and explain ideas, reasoning, and representations to one another in pair, small-group, and whole-class discourse using discipline-specific terminology, language constructs, and symbols.</li> <li>● Seek to understand the approaches used by peers by asking clarifying questions, trying out others' strategies, and describing the approaches used by others.</li> <li>● Listen carefully to and critique the reasoning of peers using data to support or counterexamples to refute arguments.</li> <li>● Develop the skills to justify mathematical reasoning with clarity and precision.</li> <li>● Practice constructing data-based arguments with specific audiences in mind.</li> <li>● Consider matters of accessibility in designing and executing their communications.</li> <li>● Consider the pros and cons of various types of data visualizations and how they fit the communicative situation.</li> </ul>	<ul style="list-style-type: none"> <li>● Introduce concepts in a way that connects students' experiences to course content and that bridges from informal contextual descriptions to formal definitions.</li> <li>● Clarify the use of data science terminology and symbols, especially those used in different contexts or different disciplines.</li> <li>● Engage students in purposeful sharing of ideas in data science, reasoning, and approaches using varied representations.</li> <li>● Support students in developing active listening skills and in asking clarifying questions to their peers in a respectful manner that deepen understanding.</li> <li>● Facilitate discourse by positioning students as authors of ideas who explain and defend their approaches.</li> <li>● Provide regular opportunities for students to communicate about data science with a variety of data visualizations</li> <li>● Scaffold instruction to support students in developing the required reading and writing skills.</li> </ul>
--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

<p><b>Technology.</b> Courses introduce students to current data science technologies and prepare them to learn and use new ones.</p>	<ul style="list-style-type: none"> <li>● Use technology to visualize and understand important data science concepts and as a tool in problem solving.</li> <li>● Understand the necessity of digital tools in cleaning and analyzing large data sets and are able to select appropriate tools for different situations.</li> <li>● Develop experience in learning new tools which allows them to try out emerging data science tools in the future.</li> <li>● Understand that the use of tools or technology does not replace the need for an understanding of reasonableness of results or how the results apply to a given context.</li> </ul>	<ul style="list-style-type: none"> <li>● Introduce students to various digital data science tools and support them in understanding the best uses for each.</li> <li>● Facilitate student learning of technological platforms through exploration, as this will aide in transferring the knowledge to future platforms.</li> <li>● Not be experts in the use of every platform but willing to experiment along with students' questions and model good practices for seeking answers to such questions.</li> </ul>
<p><b>Assessment.</b> Courses use project-based assessments to evaluate student progress.</p>	<ul style="list-style-type: none"> <li>● Assemble a collection of their work which includes both their mathematical work and reflections on their learning process and their evolving understanding of the field of data science.</li> <li>● At the end of the course, have a portfolio of data science work that showcases their knowledge of data science as well as their software skills. This</li> </ul>	<ul style="list-style-type: none"> <li>● Provide students with projects through which they are exposed to new content and demonstrate their ability to use this new content to solve problems. These will include products that demonstrate student learning both for the teacher, and to be included in the students' portfolios.</li> <li>● Evaluate student progress throughout the</li> </ul>

	<p>portfolio might be shared with a potential employer or educational institution.</p>	<p>course by considering students' evolving portfolios as well as their reflections on their learning.</p> <ul style="list-style-type: none"> <li>• In the final project of the course, allow students freedom to decide the topic and methods of their data exploration, so that they can bring together the various skills they will have developed over the course, and allow the teacher to assess their progress.</li> </ul>
--	----------------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

## 1461 Content Learning Outcomes

1462 In this section we present the mathematical content outcomes expected from a high  
 1463 school data science course. These will be motivated by realistic examples and projects  
 1464 which will help students develop their basic data science skills as well as a larger  
 1465 understanding of their contexts and of the importance of data in their lives.

### 1466 Understanding the Role of Data in The World

1467 Students demonstrate an understanding of what qualifies as data and the different types  
 1468 of data that exist. They also understand how data is generated and collected, and the  
 1469 existence of extremely large amounts of data created by our digital lives. Students  
 1470 consider their own privacy and data footprint. Throughout the course, students discuss  
 1471 the ethics and consequences of collecting and using big data, and the ways data is  
 1472 collected, including the bias that may be present in the data collection or selection  
 1473 process. Students evaluate and critique data-based claims and arguments, in particular,  
 1474 they distinguish correlation and causation. Students understand that all data and data-  
 1475 based arguments have several sources of bias and are able to identify them. They  
 1476 understand the importance of communicating with data and making data-based

1477 arguments. They use multiple different types of data visuals both for analysis and in  
1478 order to share their thinking with others.

- 1479 ● Represent data represented by real numbers using dot, box and histograms  
1480 (S&P)
- 1481 ● Summarize categorical data for two categories in two-way frequency tables  
1482 (S&P)
- 1483 ● Interpret relative frequencies in the context of the data (S&P)
- 1484 ● Recognize possible associations and trends in the data (S&P)
- 1485 ● Distinguish between correlation and causation (S&P)
- 1486 ● Evaluate the purpose of and differences between sample surveys, experiments  
1487 and observational studies and how randomization effects each (S&P)

#### 1488 Asking Data-Based Questions

1489 Students are able to identify the types of questions that are subject to exploration  
1490 through data as well as formulate their own. They are able to perform exploratory data  
1491 analyses to draw preliminary conclusions to explore further. They can do this in a  
1492 variety of platforms. Students can look at the data available and identify questions that it  
1493 can answer as well as determine what data might be collected in order to answer a  
1494 question. Students consider how they might use some of the data they have access to,  
1495 in order to predict other variables for which it might be harder to collect data directly.

#### 1496 Unraveling the Story That Data Is Telling

- 1497 ● When working with numerical data students can describe a distribution using its  
1498 shape, center, and spread. They are able to make predictions based on these  
1499 characteristics, as well as compare distributions to one another. Students are  
1500 also able to compare two numerical variables to each other using scatter plots  
1501 and can use their understanding of functions (linear, polynomial, exponential) to  
1502 fit their data to a curve (using appropriate technological tools) and use this model  
1503 to make predictions. Students are also able to work with categorical variables in  
1504 frequency tables as well as use numerical and categorical variables together in  
1505 order to answer questions about the data. Analyze the shape of data distributions

- 1506 and compare data distribution using measures of center (mean, median) and  
 1507 spread (interquartile range (IQR), standard deviation) of different data sets (S&P)
- 1508 ● Interpret differences in shape, center and spread including the effects of outliers  
 1509 (S&P)
  - 1510 ● Use mean and standard deviation to fit to a normal distribution and to estimate  
 1511 population percentages. Know that this procedure is not appropriate for all data  
 1512 sets (S&P)
  - 1513 ● Use tech tools (calculators, spreadsheets and tables) to estimate areas under the  
 1514 normal curve (S&P)
  - 1515 ● Represent two variable data on a scatter plot and describe how the variables are  
 1516 related (S&P)
  - 1517 ● Fit a linear function on scatter plots where the data suggests a linear fit (S&P)
  - 1518 ● Fit a function to the data to solve problems in context of the data (S&P)
  - 1519 ● Determine the fit of a function by plotting and analyzing residuals (S&P)
  - 1520 ● In a linear model interpret slope as a rate of change and the intercept as the  
 1521 constant term in the context of the data (S&P)
  - 1522 ● Use technology to compute and interpret the correlation coefficient of a linear fit
  - 1523 ● Estimate a line of best fit for a single linear regressions (algebra)
  - 1524 ● Determine and interpret the strength of correlation to determine the best fit.  
 1525 (algebra)
  - 1526 ● Use multiple linear regressions and non-linear regressions (linear, quadratic,  
 1527 exponential) algebra
  - 1528 ● Regression trees and classification trees (Talitha course has a lesson on these.  
 1529 “Classification trees and regression trees are easily understandable and  
 1530 transparent methods for predicting or classifying new records”)
  - 1531 ● Understand independent and dependent events and the and that two  
 1532 independent events have a probability of occurring together that is a product of  
 1533 their individual probability of occurring (S&P)
  - 1534 ● Conditional probability (S&P)
  - 1535 ● Construct and interpret two-way frequency tables of data when two categories  
 1536 are associated with each object being classified. Use the two-way table as a

- 1537 sample space to decide if events are independent and to approximate conditional  
 1538 probabilities. (S&P)
- 1539 ● Recognize and explain the concepts of conditional probability and independence  
 1540 in everyday language and situations. (S&P)
  - 1541 ● Calculate expected values and use them to solve problems. (S&P)
    - 1542 ○ Calculate the expected value of a random variable and interpret it as the  
 1543 mean of the probability distribution (S&P)
  - 1544 ● Use probability to evaluate outcomes of decisions (S&P)
  - 1545 ● Linear algebra:
  - 1546 ● Recognize situations in which one quantity changes at a constant rate per unit  
 1547 interval relative to another (algebra)
  - 1548 ● Recognize situations in which a quantity grows or decays by a constant percent  
 1549 rate per unit interval relative to another (algebra)

1550 **Grappling with Variability and Uncertainty**

1551 Students understand variability is inherent to data and are able to identify multiple  
 1552 sources of it. They practice collecting and organizing data about their own lives and  
 1553 communities as well as working with large, real-world, publicly available data sets.  
 1554 Students consider sampling practices and how they affect the data that is collected.  
 1555 They can use probability to make decisions and understand the uncertainty that comes  
 1556 along with predictions.

- 1557 ● Know that statistics is a process for making inferences about population  
 1558 parameters based on random samples of the population (S&P)
- 1559 ● Determine if a model from a data generating process or simulation is accurate  
 1560 (S&P)
- 1561 ● Make inferences and justify conclusions from sample surveys, experiments and  
 1562 observational studies (S&P)
- 1563 ● Use data from a sample to estimate population mean or proportion and develop a  
 1564 margin of error through simulation models (S&P)
- 1565 ● Use simulations to decide if differences between parameters are significant  
 1566 (S&P)

- 1567       • Evaluate reports based on data (S&P)

1568   Transforming Data with Technology

1569   Students understand that data is not always collected/shared/received ready to be  
1570   analyzed and it sometimes requires work to prepare it. They can use different digital  
1571   tools to clean and transform the data (e.g. merge data, deal with incomplete data,  
1572   normalize data). They are familiar with the basics of programming as needed, and are  
1573   comfortable editing code or finding the appropriate tools to transform the data in ways  
1574   useful to their own data analysis. Students can combine their knowledge of probability  
1575   and programming to construct simulations of probabilistic events, and they understand  
1576   the basic idea behind machine learning as well as its power and shortcomings.

- 1577       • Perform operations on matrices and use matrices in applications (Number)
- 1578       • Use matrices to represent and manipulate data (Number)
- 1579       • Cleaning names, categories, and strings (IDS)
- 1580       • Simulation using experimental data (IDS)
- 1581       • Translate between different bit representations of real-world phenomena, such as  
1582       characters, numbers, and images (CS Stds)
- 1583       • Evaluate the tradeoffs in how data elements are organized and where data is  
1584       stored.(CS Stds)
- 1585       • Create clearly named variables that represent different data types and perform  
1586       operations on their values. (CS Stds)
- 1587       • Collect data using computational tools and transform the data to make it more  
1588       useful and reliable. (CS Stds)

1589   Sample Courses

1590   Effective Data Science courses consider how to help students:

- 1591       • understand how data are used by professionals to address real-world problems;
- 1592       • understand that data are used in all facets of modern life;
- 1593       • understand how data support science to identify and tackle real-world problems  
1594       in our communities;



- 1595 • analyze statistical graphics to identify patterns in data and to connect these  
1596 patterns back to the real world;
- 1597 • understand that by treating photos, words, numbers, and sounds as data, we can  
1598 gain insight into the real world;
- 1599 • learn to analyze data, including: posing questions that can be answered by  
1600 considering relations among variables in a data set, using collected data to  
1601 generate hypotheses for future data collection, critically evaluating shortcomings  
1602 and strengths in the data and the data collection process, and informally  
1603 evaluating hypotheses using data at hand.

1604 A sample table of contents for the course is given in Appendix A.

1605 Another sample course begins with a consideration of the meaning of data, the  
1606 importance of communicating data visually, investigating community issues, cleaning  
1607 data, exploratory data analysis, ethical issues around data, creating data dashboards,  
1608 linear and nonlinear regression models, statistics, probability, and forecasting. The  
1609 course is designed to engage students actively and to be flexible enough for teachers to  
1610 include local issues of importance to their communities. While addressing concepts of  
1611 Data Analysis with rigor, the access and dependence upon current, local, and publicly  
1612 accessible data is a key feature. One goal of the course is that it be open to all students,  
1613 regardless of prior mathematics achievement, all lessons will be “low floor and high  
1614 ceiling” – designed so that everyone can access them and they extend to high levels.

1615 Some schools have created a Data Science elective for students in grades 10–12. The  
1616 course may begin with the basics of data collection, and then teach distributions, linear  
1617 regression, probability, and statistical inference through investigation-based activities.  
1618 Course activities may include making distributions of students texting-frequency,  
1619 examining player statistics from 30 Major League Baseball teams, and analyzing the  
1620 link between poverty and obesity. Districts can design their course to meet A–G  
1621 Mathematics credit requirements.

1622 An additional example of school-created course for students in grade nine is one  
1623 focused on software design and data science. It teaches algebraic, geometric, and

1624 statistical concepts through contexts like video-game design. This course can be an  
1625 example of a modernized integrated pathway, teaching the traditional sequence through  
1626 modern mediums and applications. The course can also be designed to meet A–G  
1627 elective credit requirements.

1628 The different examples of courses and high school approaches above use different  
1629 software and tools, which seems appropriate as data science does not require any  
1630 particular software package, it is more important that students learn to ask good  
1631 questions and apply an effective tool to help them answer them. Exposure to some  
1632 software is essential for those wishing to pursue a full-time career in data science, and  
1633 comfort with such programs is increasingly valuable for many other professions that  
1634 involve basic data analysis.

1635 In total, over 70 individual high schools and 15 districts offered a data science  
1636 mathematics or elective course in California during the 2019–2020 school year that  
1637 counted for A–G credit (University of California data). That compares to just 34 high  
1638 schools and 6 districts two years before in 2017–2018. This rapid increase in course  
1639 offerings is likely an indication of both high interest in and importance of data science  
1640 content throughout the curriculum.

## 1641 High School Tools and Resources

1642 One sample tool for students to explore large data sets is the free, open source  
1643 software tool Common Online Data Analysis Platform (CODAP)  
1644 (<https://learn.concord.org/dynamic-data-science>) from the Concord Consortium. Using  
1645 this software, students can import data from their own community or work with the large  
1646 data sets already available in the tool. Students will learn to become active citizens in  
1647 their communities, learning that mathematics is an important tool for benefitting their  
1648 community.

1649 The Census at School project (<https://ww2.amstat.org/censusatschool/>) is an  
1650 international classroom project that engages students in grades 4–12 in statistical  
1651 problem solving. Students complete a brief online survey, analyze their class census

1652 results, and compare their class with random samples of students in the United States  
1653 and other countries.

1654 Other software such as Fathom (<https://fathom.concord.org/>) and Statkey  
1655 (<http://www.lock5stat.com/StatKey/>) allow exploration and organization of data sets and  
1656 the development of simulations. Google offers a free coding software called Google  
1657 Script Coding, Python is another tool that can be used to explore and analyze data sets.  
1658 A more sophisticated data software tool is R. This requires learning time and schools  
1659 may need to provide server space to run the software.

1660 Below are two examples of data science projects that students may work on in high  
1661 school, freely available from the Concord Consortium:

1662 In the California American Community Survey (ACS) Data Portal  
1663 (<https://learn.concord.org/dynamic-data-science>) students are given access to the data  
1664 portal which gives census data for California residents from the U.S. Census Bureau's  
1665 American Community Survey. The database contains demographic information about  
1666 California residents (e.g., marital status, sex, place of birth, employment status, and  
1667 health information). Data challenges are given such as finding out the average income  
1668 of Californians of different age groups in 2013, or students can choose to investigate  
1669 their own questions. For example, they may choose to look at salaries by gender, or  
1670 make a data visualization to show the different ethnic groups that live in California.  
1671 Standards addressed include making inferences and justifying conclusions (HSS-  
1672 IC.A.1) and SMP.2 (Reason abstractly and quantitatively), SMP.3 (Construct viable  
1673 arguments and critique the reasoning of others), SMP.4 (Model with mathematics), and  
1674 SMP.5 (Use appropriate tools).

1675 In a different Concord activity: Making Trees in a Diagnosis Game  
1676 (<https://learn.concord.org/resources/1241/trees-in-a-diagnosis-game>) students use data  
1677 to build binary trees for decision-making and prediction. Prediction trees are the first  
1678 step towards linear regression, which plays an important role in machine learning for  
1679 future data scientists. Students begin by manually putting "training data" through an  
1680 algorithm. They then learn to automate the process and to test their ability to predict

1681 which alien creatures are sick and which are healthy. This activity touches upon many  
1682 content and practice standards, including Making inferences and justifying conclusions  
1683 (HSS.-IC.A.1), Using probability to make decisions (HSS.-MD.B.7), and all standards for  
1684 mathematical practice.

## 1685 Conclusion

1686 Changing demands for life in a data-rich world require that California schools prepare all  
1687 students to examine claims justified with data, to understand the probabilistic  
1688 underpinning of drawing conclusions from samples, and to see data as a tool to answer  
1689 many questions of interest. Developing these abilities requires that students generate  
1690 questions and work with data beginning in kindergarten (or before), and have  
1691 experiences of increasing depth and complexity throughout their school careers.  
1692 Students who wish to focus extra attention on data science should have an opportunity  
1693 to pursue advanced courses late in their high school careers.

1694 Above all, students at all levels should have experiences that build their mathematical  
1695 toolkit for making sense of their worlds.

## Free Resources for the Teaching of Data Science

- Concord Consortium: <https://learn.concord.org/dynamic-data-science>
- Jo Boaler Online Course: The teaching of data science K–12:  
<https://www.youcubed.org/21st-century-teaching-and-learning/>
- The Messy Data Coalition: <https://messydata.org/>
- University of Chicago RISC: <https://www.21cmath.org/>
- Women in Data Science Video:  
<https://www.youcubed.org/resources/what-is-data-science/>
- Wolfram-Alpha: <http://www.computerbasedmath.org/>
- Youcubed Resources: <https://www.youcubed.org/resource/data-literacy/>
- Youcubed Grades 6–10 Data Lessons:  
<https://www.youcubed.org/data-science-lessons/>
- Youcubed Data Talks:  
<https://www.youcubed.org/resource/data-talks/>

## 1696 References

- 1697 Arnold, P. (2007). What about the P in the PPDAC cycle? An initial look at posing  
1698 questions for statistical investigation. *Education*, 55.
- 1699 Bargagliotti, A., Franklin, C., Arnold, P., Gould, R., Johnson, S., Perez, L., Spangler, D.  
1700 (2020). *Pre-K–12 Guidelines for Assessment and Instruction in Statistics Education II*  
1701 *(GAISE II): A Framework for Statistics and Data Science Education*. American  
1702 Statistical Association.
- 1703 Boaler (2019). *Limitless Mind. Learn, Lead and Live without Barriers*. Harper Collins.
- 1704 Boaler, J., Cordero, M., & Dieckmann, J. (2019). Pursuing Gender Equity in  
1705 Mathematics Competitions. A Case of Mathematical Freedom. Mathematics Association

1706 of America, FOCUS, Feb/March 2019.  
1707 [http://digitaleditions.walsworthprintgroup.com/publication/?m=7656&l=1#{%22issue\\_id%](http://digitaleditions.walsworthprintgroup.com/publication/?m=7656&l=1#{%22issue_id%22:566588,%22page%22:18)  
1708 [22:566588,%22page%22:18](http://digitaleditions.walsworthprintgroup.com/publication/?m=7656&l=1#{%22issue_id%22:566588,%22page%22:18)

1709 Carmichael, I., Marron, J.S. Data science vs. statistics: two cultures?. *Jpn J Stat Data*  
1710 *Sci* 1, 117–138 (2018). <https://doi.org/10.1007/s42081-018-0009-3>

1711 Chestnut, E. K., Lei, R. F., Leslie, S. J., & Cimpian, A. (2018). The myth that only  
1712 brilliant people are good at math and its implications for diversity. *Education sciences*,  
1713 8(2), 65.

1714 CORE SEL Competencies: <https://casel.org/core-competencies/>

1715 Dhar, V. (2013). Data science and prediction. *Communications of the ACM*, 56(12), 64-  
1716 73. [https://cacm.acm.org/magazines/2013/12/169933-data-science-and-](https://cacm.acm.org/magazines/2013/12/169933-data-science-and-prediction/fulltext)  
1717 [prediction/fulltext](https://cacm.acm.org/magazines/2013/12/169933-data-science-and-prediction/fulltext)

1718 D'Ignazio, C., & Klein, L. F. (2020). *Data feminism*. MIT Press.

1719 Paris, D. (2012). Culturally sustaining pedagogy: A needed change in stance,  
1720 terminology, and practice. *Educational researcher*, 41(3), 93-97.

1721 Pelesko, John (2015). “‘The’ Modeling Cycle.”  
1722 <http://modelwithmathematics.com/2015/08/the-modeling-cycle/>

1723 Rubin, Andee. "Learning to Reason with Data: How Did We Get Here and What Do We  
1724 Know?." *Journal of the Learning Sciences* (2019): 1–11

1725 Walton, G. M., Logel, C., Peach, J. M., Spencer, S. J., & Zanna, M. P. (2015). Two brief  
1726 interventions to mitigate a “chilly climate” transform women’s experience, relationships,  
1727 and achievement in engineering. *Journal of Educational Psychology*, 107(2), 468.